



Mean Reversion in Housing Markets

Citation

Nathanson, Charles Gordon. 2014. Mean Reversion in Housing Markets. Doctoral dissertation, Harvard University.

Permanent link

<http://nrs.harvard.edu/urn-3:HUL.InstRepos:12274618>

Terms of Use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Other Posted Material, as set forth at <http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#LAA>

Share Your Story

The Harvard community has made this article openly available.
Please share how this access benefits you. [Submit a story](#).

[Accessibility](#)

Mean Reversion in Housing Markets

A dissertation presented

by

Charles Gordon Nathanson

to

The Department of Economics

in partial fulfillment of the requirements

for the degree of

Doctor of Philosophy

in the subject of

Economics

Harvard University

Cambridge, Massachusetts

April 2014

© 2014 Charles Gordon Nathanson

All rights reserved.

Dissertation Advisor:

Professor Edward L. Glaeser

Author:

Charles Gordon Nathanson

Mean Reversion in Housing Markets

Abstract

Booms in house prices are usually followed by busts. This pattern is called “mean reversion.” Mean reversion in housing markets has historically coincided with economic recessions across the world. Chapter 1 establishes mean reversion in U.S. data, and attempts to explain it using the dynamics of wages in cities. Chapter 2 takes a different approach. It models mean reversion resulting from speculation and uncertainty. This model explains why strong mean reversion in prices occurs in cities where it is easy to build houses, a phenomenon that Chapter 1 cannot explain. Chapter 3 takes the spirit of Chapter 2 and applies it to the optimal design of the income tax.

Contents

Abstract	iii
Acknowledgments	ix
1 Housing Dynamics: An Urban Approach	1
1.1 Introduction	1
1.2 A Dynamic Model of Housing Prices	5
1.2.1 Housing Supply	5
1.2.2 Housing Demand	6
1.2.3 Equilibrium	7
1.3 Estimating the Model	12
1.3.1 Data	16
1.3.2 Methodology	17
1.3.3 Estimation Results	20
1.4 Matching the Data and Discussion	22
1.4.1 The Impact of Information on the Predictions of the Model	22
1.4.2 Volatility and Serial Correlation in House Prices	24
1.4.3 Volatility and Serial Correlation in Construction	28
1.5 Conclusion	29
2 Arrested Development: Theory and Evidence of Supply-Side Speculation in the Housing Market	31
2.1 Introduction	31
2.2 Stylized Facts of the U.S. Housing Boom and Bust	38
2.2.1 The Cross-Section of Cities	38
2.2.2 The Central Importance of Land Prices	41
2.2.3 Land Market Speculation by Homebuilders	42
2.3 A Housing Market with Homeowners and Developers	46
2.4 Supply-Side Speculation	52
2.4.1 Land Speculation and Dispersed Homeownership	53
2.4.2 Belief Aggregation	55
2.5 The Cross-Section of City Experiences During the Boom	58

2.6	Variation in House Price Booms Within Cities	65
2.6.1	Location	65
2.6.2	Structure Type	66
2.7	Conclusion	67
3	Taxation and the Allocation of Talent	69
3.1	Introduction	69
3.2	Theory	73
3.2.1	General model	74
3.2.2	A simple case	79
3.3	Calibration	82
3.3.1	Income distributions	84
3.3.2	Externality shares	89
3.3.3	Results	94
3.4	Discussion	97
3.4.1	Allocation of talent	97
3.4.2	Labor supply elasticity debate	98
3.4.3	Debates on taxation outside neoclassical economics	100
3.4.4	Closely related literature	101
3.5	Structural Model with General Ability	103
3.5.1	Calibration	104
3.5.2	Optimal tax rates	108
3.5.3	Quantitative importance of elasticities vs. externalities	111
3.5.4	Quantitative welfare gains	113
3.5.5	Effects of the Reagan tax reforms	115
3.6	Conclusion	118
	References	120
	Appendix A Appendix to Chapter 1	129
A.1	Estimation Details	129
A.1.1	Sequential Two-Step GMM Estimator	129
A.1.2	Moment Conditions	131
A.1.3	Stochastic Processes Predicted by the Model	132
A.2	Definitions of Trend Variables	133
A.3	Proofs	134
A.3.1	Proof of Lemma 1	134
A.3.2	Proof of Lemma 2	136
A.3.3	Proof of Proposition 1	138

A.3.4	Proof of Proposition 2	138
A.3.5	Proof of Proposition 3	139
A.4	Calculation of Volatilities in Table 1.1	140
A.5	BEA Income Data Tables	141
Appendix B	Appendix to Chapter 2	144
B.1	Micro-foundation of owner-occupancy utility	144
B.2	Proofs	145
B.2.1	Proof of Lemma 3	145
B.2.2	Proof of Proposition 5	147
B.2.3	Proof of Proposition 6	148
B.2.4	Proof of Implication 7	151
B.2.5	Proof of Implication 8	151
B.2.6	Proof of Implication 9	151
B.3	Construction equation	152
Appendix C	Appendix to Chapter 3	153
C.1	Alternative Finance Calibration	153
C.2	Externality Share Calibration	156
C.2.1	Law and Computers/Engineering	156
C.2.2	Management	158
C.2.3	Academia/Science	159
C.2.4	Consulting	160
C.2.5	Teaching	161
C.2.6	Arts/Entertainment	162
C.3	General Ability Model	162
C.4	Estimation	165
C.5	Solving for the optimal tax function	170
C.6	Alternative Elasticity Value	173
C.7	Allocation of Talent in the General Ability Model	174
C.8	Proofs	176

List of Tables

1.1	Relative Volatility of Terms in House Price Equation	15
1.2	Estimated Demand and Supply Parameters HMDA Income Data, 1990- 2004	20
1.3	Sensitivity of Predictions to Different Information Structures	24
1.4	Volatility and Serial Correlation in House Prices and Construction: HMDA Income Data, 1990-2004	25
3.1	Sources of externality estimates from the economics literature	89
3.2	Externality profiles in each of four calibrations	91
3.3	Welfare Gains	114
3.4	Reallocation of Talent from Reagan Tax Reforms	116
A.1	Estimated Demand and Supply Parameters: BEA Income Data, 1980-2003 . .	142
A.2	Volatility and Serial Correlation in House Prices and Construction: BEA Income Data, 1980-2003	143
C.1	Real net income in finance over time	167
C.2	Welfare Gains Under Alternate Elasticity Value	173

List of Figures

1.1	Real House Price Appreciation in the 1980s and 1990s	3
1.2	Housing Unit Growth in the 1980s and 1990s	4
1.3	Simulated One-Time Shock	11
2.1	Long-Run Development Constraints in Las Vegas	36
2.2	The U.S. Housing Boom and Bust Across Cities	39
2.3	Input Price and House Price Increases Across Cities, 2000-2006	43
2.4	Supply-Side Speculation Among U.S. Public Homebuilders, 2001-2010	45
2.5	Model Simulations For Different Cities	61
3.1	Income distributions fitted to IRS and Harvard data in 9 industries	86
3.2	Baseline and alternative income calibrations in Finance	87
3.3	The allocation of talent condition on income level	88
3.4	Social Product in Different Professions	93
3.5	ATEM and MTEM marginal tax rates	95
3.6	Reference Income Distributions	107
3.7	Optimal Tax Rates in Structural Model	109
3.8	Horserace between elasticities and externalities	116
C.1	Income Distributions Under Alternative Finance Calibrations	154
C.2	ATEM and MTEM Policies Under Alternate Finance Calibration	155
C.3	“Horserace” Under Alternate Finance Calibration	155
C.4	Optimal Tax Rates Under Different Elasticity Calculation	172
C.5	Allocation of talent by ability quantile under <i>laissez-faire</i>	175
C.6	Allocation of Talent in Different Regimes	175
C.7	Reallocation of Talent from Reagan Tax Reforms	176

Acknowledgments

I thank my coauthors Ed Glaeser, Joe Gyourko, Ben Lockwood, Eduardo Morales, Glen Weyl, and Eric Zwick for writing this dissertation with me. My committee—Professors John Campbell, Ed Glaeser, David Laibson, and Andrei Shleifer—provided invaluable guidance, advice, and support for which I am grateful. Several other faculty gave helpful feedback on this work: Professors Raj Chetty, Claudia Goldin (who provided data used in Chapter 3), Robin Greenwood, Sam Hanson, Alp Simsek, Adi Sunderam, and Jeremy Stein. Finally, I acknowledge financial support from the NSF Graduate Research Fellowship Program, the Bradley Foundation, the Alfred P. Sloan foundation, and the Becker Friedman Institute at the University of Chicago.

Chapter 1

Housing Dynamics: An Urban Approach¹

1.1 Introduction

Can the dynamics of housing markets be explained by a dynamic, rational expectations version of the standard urban real estate models of Alonso (1964), Rosen (1979), and Roback (1982)? In this tradition, housing prices reflect a spatial equilibrium, where prices are determined by local wages and amenities so that local heterogeneity is natural. Our model extends the Alonso-Rosen-Roback framework by focusing on high frequency price dynamics and by incorporating endogenous housing supply.

An urban approach can potentially help address the fact that most variation in housing price changes is local, not national. Less than eight percent of the variation in price levels and barely more than one-quarter of the variation in price changes across cities can be accounted for by national, year-specific fixed effects. Clearly, there is much local variation that cannot be accounted for by common macroeconomic variables such as interest rates or national income.

We focus not on the most recent boom and bust, which was extraordinary in many

¹This chapter is co-authored with Edward L. Glaeser, Joseph Gyourko, and Eduardo Morales.

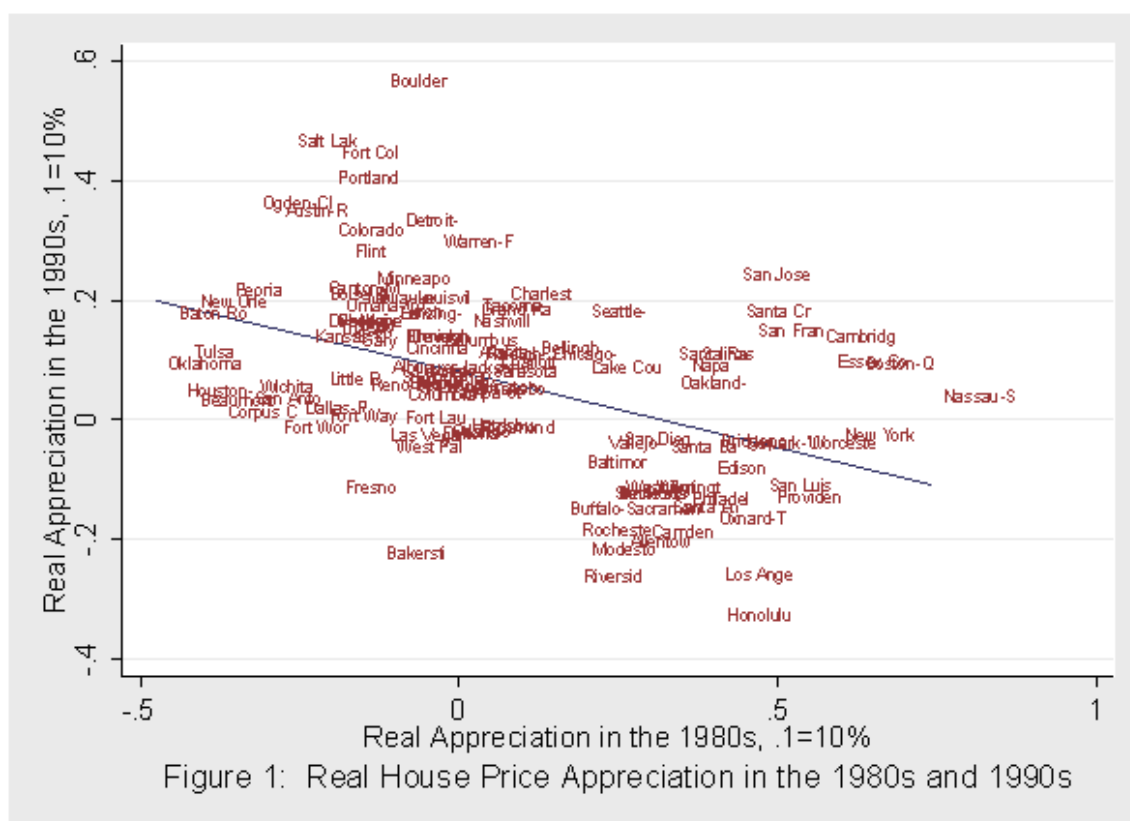
dimensions, but rather on long-term stylized facts about housing markets. One such fact is that price changes are predictable (Case and Shiller, 1989; Cutler *et al.*, 1991). Depending upon the market and specific time period being examined, a \$1 increase in real constant quality house prices in one year is associated with a 60-80 cent increase the next year. However, a \$1 increase in local market prices over the past five years is associated with strong mean reversion over the next five year period. This raises the question of whether the high frequency momentum and low frequency mean reversion of price changes can be reconciled with a rational market.

Another outstanding feature of housing markets is that the strong mean reversion in price appreciation and strong persistence in housing unit growth across decades shown in Figures 1.1 and 1.2 is at odds with simple demand-driven models in which prices and quantities move symmetrically. This raises the question of what else is needed to generate this pattern.

Third, price changes and construction levels are quite volatile in many markets. The range of standard deviations of three-year real changes in our sample of metropolitan area average house prices runs from about \$6,500 in sunbelt markets to over \$30,000 in coastal markets. New construction within markets also can be volatile, with its standard deviation much higher in the sunbelt region. Can this volatility be the result of real shocks to housing markets or must it reflect bubbles or animal spirits?

Section 2 presents our model and its implications. Naturally, the urban approach predicts that housing markets are local, not national, in nature. Predictable housing price changes also are shown to be compatible with a no-arbitrage rational expectations equilibrium. Mean reversion over the medium and longer term results if construction does not respond immediately to shocks and if local income shocks themselves mean revert. High frequency positive serial correlation of housing prices results if there is enough positive serial correlation of labor demand or amenity shocks. Conceptually, a dynamic rational expectations urban model is at least consistent with the outstanding features of housing markets, at least as they existed prior to the financial crisis.

Figure 1.1: Real House Price Appreciation in the 1980s and 1990s



However, our calibration exercises yield both successes and failures in trying to match key moments of the data. We are able to capture the extensive heterogeneity across different types of markets, especially in our contrast of coastal markets with high inelastic supply sides with interior markets with very elastic supplies of homes. Different shocks to the varying local income processes interact with very different supply side conditions to generate materially different housing market dynamics.

The model also does a reasonably good job of generating high variation in house price changes based on innovations in our proxy for local incomes, although we cannot match the extremely high volatility in house prices in the most variable coastal markets. The model also does a tolerably good job of matching the volatility of new construction, generating wide divergences across markets based on underlying supply elasticities. However, the model again cannot match the most volatile construction markets which are off the coasts.

Figure 1.2: *Housing Unit Growth in the 1980s and 1990s*

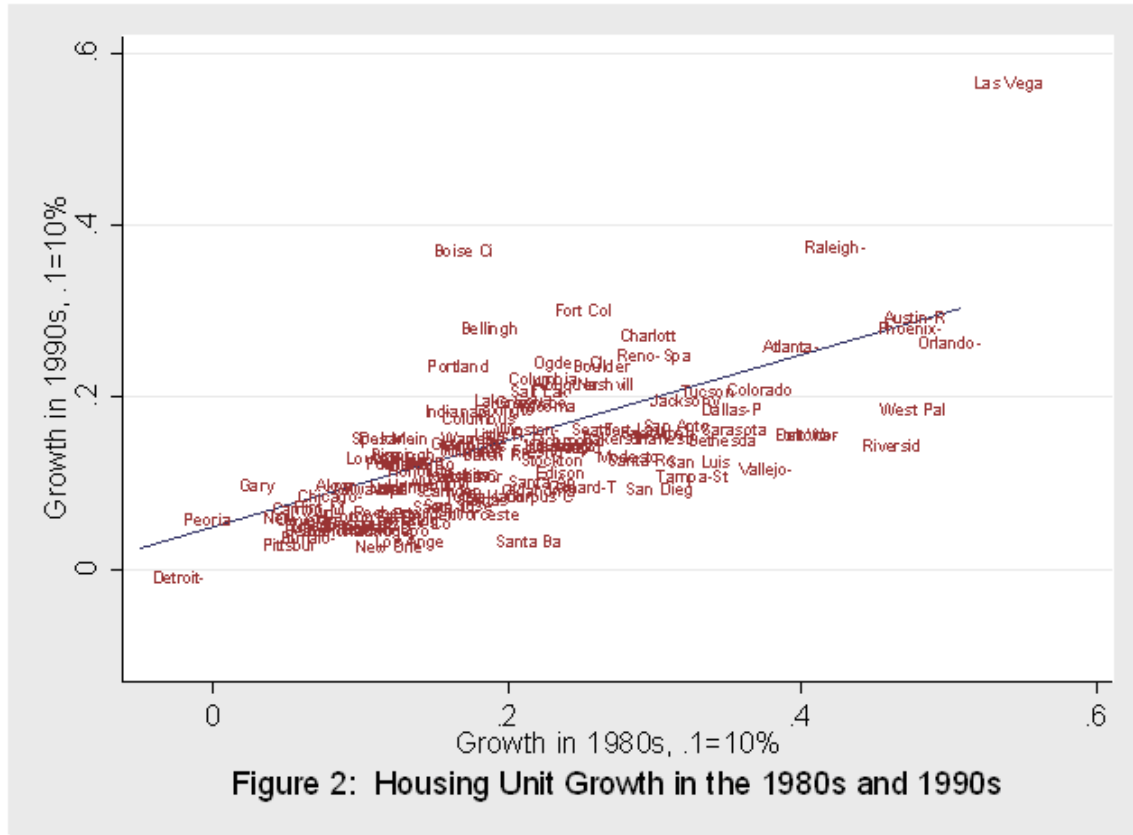


Figure 2: *Housing Unit Growth in the 1980s and 1990s*

With respect to the serial correlations of quantities and prices, the model gets the pattern, but not the magnitude, of the strong high-frequency persistence in construction. Our model correctly captures the weakening of that persistence over longer horizons, but still cannot replicate the mean reversion that is evident in the data over five-year periods. The model fails utterly at explaining the very strong, high frequency positive serial correlation in price changes. It does a better job at predicting mean reversion over longer five-year horizons, but still cannot precisely match the magnitude of that pattern, especially in coastal markets.

This suggests that the most important puzzle for housing economists to explain, apart from the most recent cycle, is the strong persistence in high frequency price changes from one year to the next. Persistence itself is not enough to reject a rational expectations model, but the mismatch between the data and model at annual frequencies indicates that Case

and Shiller (1989)'s conclusion regarding inefficiency could be right. Other issues deserving closer examination include whether there really is excess volatility in coastal markets and the nature of serial correlation in construction over longer time horizons.

1.2 A Dynamic Model of Housing Prices

1.2.1 Housing Supply

Homebuilders are risk neutral firms that operate in a competitive market. Suppressing a subscript for individual markets for ease of exposition, the marginal cost to this industry of constructing a house at time t is given by

$$C + c_0t + c_1I_t + c_2N_t,$$

where I_t is the amount of construction and N_t is the housing stock at time t . The c_0 term allows unit costs to trend over time. When $c_1 > 0$, the supply curve at time t is upward-sloping. The coefficient c_2 allows unit costs to depend on the city size, reflecting community opposition to development as density levels increase. We assume that $c_1 > c_2$ so that present construction has a larger effect on costs through the first effect. The supply parameters c_0 , c_1 , and c_2 can vary across metropolitan areas.

Housing is completely durable, and new supply is constrained to be non-negative:

$$I_t \geq 0.$$

Homebuilders also face a time to build. Housing constructed at time t cannot be sold until time $t + 1$. Homebuilders also bear the costs of time t construction at time $t + 1$. Perfect competition and risk-neutrality deliver the following supply condition:

$$E(H_{t+1}) = C + c_0t + c_1I_t + c_2N_t \tag{1.1}$$

when $I_t > 0$, where H_{t+1} is the house price at time $t + 1$. In equilibrium, the expected sales price of a house equals the marginal cost when homebuilders construct new houses.

1.2.2 Housing Demand

Each person consumes exactly one unit of housing, so that N_t equals both the housing stock and the population. Consumer utility depends linearly on consumption and city-specific amenities:

$$U(\text{Consumption}_t, \text{Amenities}_t) = \text{Consumption}_t + \text{Amenities}_t.$$

Consumers are identical and face a city-specific labor demand curve of

$$\text{Wages}_t = W_t - \alpha_W N_t$$

at time t . Amenities also depend linearly on the population:

$$\text{Amenities}_t = A_t - \alpha_A N_t.$$

Consumers must own a house to access the city's labor market and amenities. We exclude rental contracts from the model to focus on the owner-occupancy market. Consumers are risk-neutral and can borrow and lend at an interest rate r . Their indirect utility is therefore

$$V_t = W_t + A_t - (\alpha_W + \alpha_A)N_t - \left(H_t - \frac{E(H_{t+1})}{1+r} \right). \quad (1.2)$$

To pin down this utility level, we turn to the cross-metropolitan area spatial equilibrium concept introduced by Rosen (1979) and Roback (1982). Consumers are indifferent across cities at all points in time. This indifference condition is a particularly strong version of the standard spatial equilibrium assumption that assumes away moving costs. There is a “reservation” city where housing is completely elastic: $c_0 = c_1 = c_2 = 0$, so that housing prices always equal C .² Wages and amenities do not depend on the reservation city population: $\alpha_W = \alpha_A = 0$. If we let \bar{V}_t equal $W_t + A_t$ for the reservation city, then the

²While it is possible that prices will deviate around this value because of temporary over- or under-building, we simplify and assume that the price of a house always equals C .

reservation utility level that holds in this city as well as in all cities is

$$V_t = \bar{V}_t - \frac{rC}{1+r}. \quad (1.3)$$

The existence of the reservation city makes our calculations considerably easier, and there are places within the United States, especially in the growing areas of the sunbelt, that are marked by elastic labor demand and housing supply Glaeser *et al.* (2005).³

Putting together equations (1.2) and (1.3) gives the following housing demand equation:

$$H_t - \frac{E(H_{t+1})}{1+r} - \frac{rC}{1+r} = W_t + A_t - (\alpha_W + \alpha_A)N_t - \bar{V}_t. \quad (1.4)$$

For our estimation, we assume the following functional form:

$$W_t + A_t - \bar{V}_t = \bar{x} + qt + x_t,$$

where x_t is a stochastic term that follows an ARMA(1,1) process:

$$x_t = \delta x_{t-1} + \epsilon_t + \theta \epsilon_{t-1},$$

with $0 < \delta < 1$ and the ϵ_t independently and identically distributed with mean 0 and finite variance σ^2 . The \bar{x} term is a city fixed effect and q is a city-specific drift term. We also define

$$\alpha \equiv \alpha_W + \alpha_A$$

to be the slope of the housing demand curve, and we assume that $\alpha > 0$.

1.2.3 Equilibrium

The supply equation (1.1) and the demand equation (1.4) jointly determine equilibrium prices, housing stock, and investment. To obtain a unique solution to our model, we impose a transversality condition

$$\lim_{j \rightarrow \infty} \frac{E_t(H_{t+j})}{(1+r)^j} = 0 \quad (1.5)$$

³Van Nieuwerburgh and Weill (2010) present a similar model in their exploration of long run changes in the distribution of income.

for all t . The transversality condition limits the possible role of housing bubbles in accounting for housing dynamics. While we do not discount the possible explanatory power of bubbles, our focus here allows us to learn what aspects of housing dynamics can already be explained by a model in which prices equal the discounted sum of current and future expected rents. The following lemma shows that price, housing stock, and investment converge towards “trend” levels of these variables when the transversality condition is satisfied.

Lemma 1. *When equation (1.5) is satisfied, there exist unique price, stock, and investment functions \hat{H}_t , \hat{N}_t , and \hat{I}_t such that*

$$\lim_{j \rightarrow \infty} E_t(H_{t+j}) - \hat{H}_{t+j} = \lim_{j \rightarrow \infty} E_t(N_{t+j}) - \hat{N}_{t+j} = \lim_{j \rightarrow \infty} E_t(I_{t+j}) - \hat{I}_{t+j} = 0$$

for any H_t , N_t , and I_t that satisfy the supply and demand equations. \hat{H}_t and \hat{N}_t are linear in t and \hat{I}_t is constant.

We call \hat{H}_t , \hat{N}_t , and \hat{I}_t trend prices, population, and investment. Closed-form expressions for these trend variables as well as a proof of the lemma appear in the technical appendix.

If $x_t = 0$ for all t and $N_t = \hat{N}_t$ for some initial period, then the steady state quantities would fully describe the equilibrium.⁴ Secular trends in housing prices come from the trend in housing demand as long as $c_2 > 0$, or the trend in construction costs as long as $\alpha > 0$. If $c_2 = 0$, so that construction costs don’t increase with total city size, then trends in wages or amenities will impact city size but not housing prices. If $\alpha = 0$ and city size doesn’t decrease wages or amenities, then trends in construction costs will impact city size but not prices.

Lemma 2 then describes housing prices and investment when there are shocks to demand and when $N_t \neq \hat{N}_t$. The proof is in the technical appendix.

Lemma 2. *At time t , housing prices equal*

$$H_t = \hat{H}_t + x_t + \frac{E_t(x_{t+1})}{\bar{\phi} - \delta} - \frac{\alpha(1+r)}{1+r-\phi}(N_t - \hat{N}_t)$$

⁴In this case, the assumption that there is always some construction requires that $q(1+r) > rc_0$.

and investment equals

$$I_t = \hat{I} + \frac{1+r}{c_1} \frac{E_t(x_{t+1})}{\bar{\phi} - \delta} - (1 - \phi)(N_t - \hat{N}_t)$$

where $\bar{\phi} > 1 > \phi > 0$ are parameters that depend on α , c_1 , c_2 , and r .⁵

This lemma describes the movement of housing prices and construction around their trend levels. A temporary shock, ϵ , will increase housing prices by $(\bar{\phi} + \theta)/(\bar{\phi} - \delta)$ and increase construction by $(1+r)(\delta + \theta)/(c_1(\bar{\phi} - \delta))$. Higher values of δ (i.e., more permanent shocks) will make both of these effects stronger. Higher values of c_1 mute the construction response to shocks and increase the price response to a temporary shock (by reducing the quantity response). These results provide the intuition that places which are quantity constrained should have less construction volatility and more price volatility.

The following proposition provides implications about expected housing price changes.

Proposition 1. *At time t , the expected home price change between time t and $t + j$ is*

$$\begin{aligned} \hat{H}_{t+j} - \hat{H}_t + \frac{E_t(x_{t+1})}{\bar{\phi} - \delta} & \left(\frac{1+r}{c_1} \frac{\delta^{j-1}((1-\delta)c_1 - c_2) - \phi^{j-1}((1-\phi)c_1 - c_2)}{\phi - \delta} - 1 \right) \\ & - x_t + \left(\frac{\alpha(1+r)}{1+r-\phi} - \phi^{j-1}((1-\phi)c_1 - c_2) \right) (N_t - \hat{N}_t), \end{aligned}$$

the expected change in the city housing stock between time t and time $t + j$ is

$$j\hat{I} + \frac{1+r}{c_1(\bar{\phi} - \delta)} \frac{\phi^j - \delta^j}{\phi - \delta} E_t(x_{t+1}) - (1 - \phi^j)(N_t - \hat{N}_t),$$

and expected time $t + j$ construction is

$$\hat{I} + \frac{1+r}{c_1(\bar{\phi} - \delta)} \left(\frac{\delta^{j-1}(1-\delta) - \phi^{j-1}(1-\phi)}{\phi - \delta} \right) E_t(x_{t+1}) - \phi^{j-1}(1 - \phi^j)(N_t - \hat{N}_t).$$

Proposition 1 delivers the implication that a rational expectations model of housing

⁵The formulas for $\bar{\phi}$ and ϕ are

$$\begin{aligned} \bar{\phi} &= \frac{(1+r)(\alpha + c_1) + c_1 - c_2 + \sqrt{((1+r)(\alpha + c_1) + c_1 - c_2)^2 - 4(1+r)c_1(c_1 - c_2)}}{2c_1}, \\ \phi &= \frac{(1+r)(\alpha + c_1) + c_1 - c_2 - \sqrt{((1+r)(\alpha + c_1) + c_1 - c_2)^2 - 4(1+r)c_1(c_1 - c_2)}}{2c_1}. \end{aligned}$$

prices is fully compatible with predictability in housing prices. If utility flows in a city are high today and expected to be low in the future, then housing prices will also be expected to decline over time. Any predictability of wages and construction means that predictability in housing price changes will result in our model.

The predictability of construction and prices comes in part from the convergence to trend values. If $x_t = \epsilon_t = 0$ and initial population is above its trend level, then prices and investment are expected to converge on their trend levels from above. If initial population is below its trend level and $x_t = \epsilon_t = 0$, then price and population are expected to converge to their trend levels from below. The rate of convergence is determined by r, α, c_1 , and c_2 . Higher levels of c_1 and c_2 cause the rate of convergence to slow by reducing the extent that new construction responds to changes in demand.

The impact of a one-time shock is explored in the next proposition.

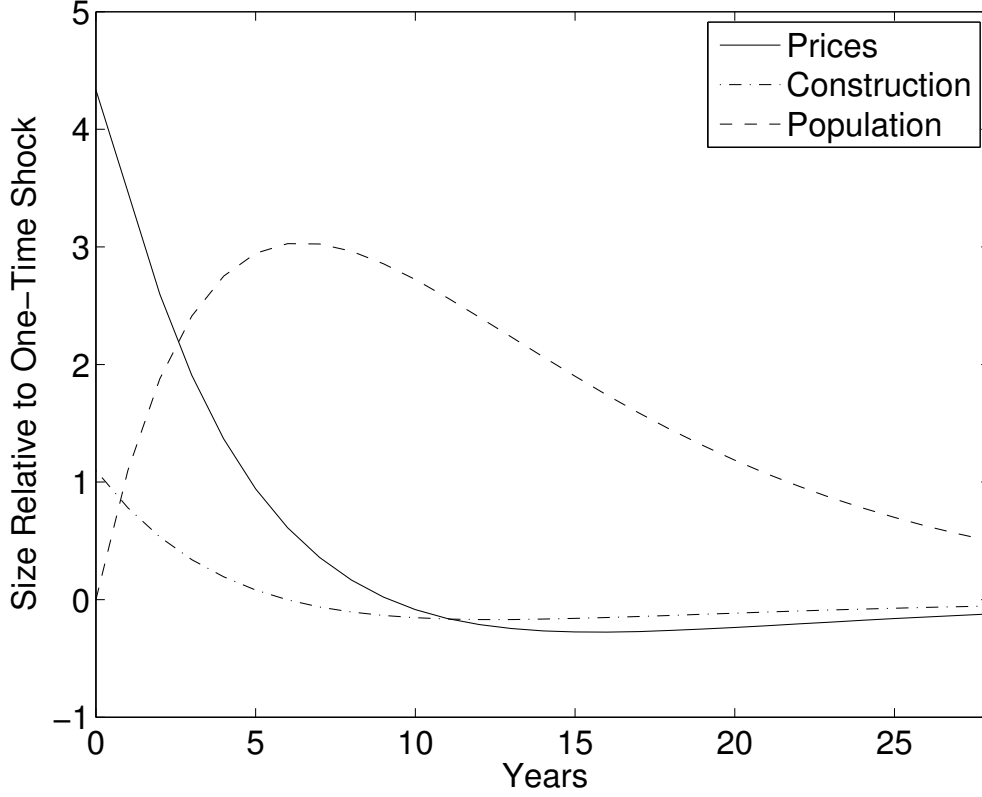
Proposition 2. *If $N_t = \hat{N}_t$, $x_{t-1} = \epsilon_{t-1} = 0$, $c_2 = 0$, and $\epsilon_t > 0$, then investment and housing prices will initially be higher than steady state levels, but there exists a value j^* such that for all $j > j^*$, time t expected values of time $t + j$ construction and housing prices will lie below steady state levels. The situation is symmetric when $\epsilon_t < 0$.*

Proposition 2 highlights that this model not only delivers mean reversion, but overshooting. Figure 1.3 shows the response of population, construction and prices relative to their steady state levels in response to a one time shock. Construction and prices immediately shoot up, but both start to decline from that point. At first, population rises slowly over time, but as the shock wears off, the heightened construction means that the city is too large relative to its steady state level. Eventually, both construction and prices end up below their steady state levels because there is too much housing in the city relative to its wages and amenities. Places with positive shocks will experience mean reversion, with a quick boom in prices and construction, followed by a bust.⁶

Finally, we turn to the puzzling empirical fact that there was strong mean reversion of

⁶Overshooting occurs here with no depreciation in the housing stock. The case with depreciation is addressed in Glaeser and Gyourko (2005).

Figure 1.3: *Simulated One-Time Shock*



Notes: We use the parameters estimated for the Interior Region using HMDA data in this figure: $\delta = 0.88$, $\theta = 0.20$, $c_1 = 3.16$, and $c_2 = 0.12$.

prices and strong positive serial correlation in population levels across the 1980s and 1990s. We address this by looking at the one period covariance of price and population changes. We focus on one period for simplicity, but we think of this proposition as relating to longer time periods. Since mean reversion dominates over long time periods, we assume $\theta = 0$ to avoid the effects of serial correlation:

Proposition 3. *If $N_0 = \hat{N}_0$, $\theta = 0$, $x_0 = \epsilon_0$, cities differ only in their demand trends q and their shock terms ϵ_0 , ϵ_1 , and ϵ_2 , and these terms are uncorrelated, then if $\delta > 1 - \phi$, second period population growth will always be positively correlated with first period population growth, while second period price growth will be negatively correlated with first period population growth as long*

as $\frac{\text{Var}(q)}{\text{Var}(\epsilon_t)}$ is below a bound.

Proposition 3 tells us that, in the model, positive serial correlation of construction levels is quite compatible with negative serial correlation of price changes. The proposition only proves that the reversal occurs when persistence of shocks is high, but in the Technical Appendix, we show that the persistence can occur when the process is less persistent. The positive correlation of quantities is driven by the heterogeneous trends in demand across urban areas. As long as the variance of these trends is high enough relative to the variance of temporary shocks, there will be positive serial correlation in quantities, as in Figure 1.2.

Yet these long trends may have little impact on price changes, since the trends are completely anticipated. As discussed above, when c_2 is low, trends will have little impact on steady state price growth, although these trends will determine the steady state price level. Instead, price changes will be driven by the temporary shocks, and if these shocks mean revert, then so will prices.

This suggests two requirements for the observed positive correlation of quantities and negative correlation of prices: city-specific trends must differ significantly and the impact of city size on construction costs must be small. Both conditions appear to occur in reality. The extensive heterogeneity in city-specific trends is discussed and documented by Gyourko *et al.* (2013) and Van Nieuwerburgh and Weill (2010). The literature on housing investment suggests that the impact of city size on construction costs is quite small (Topel and Rosen, 1988; Gyourko and Saiz, 2006).

1.3 Estimating the Model

We now calibrate the model to see whether certain moments of the data are compatible with our framework. We focus on movements in prices and construction intensity around steady state levels. The aim of this exercise is to show how a model which posits that variation in prices and construction levels is solely driven by exogenous shocks to both amenity levels and the demand for labor can fit certain moments of the housing data. As we lack data on

the short term fluctuations in the level of amenities, we will identify the parameters of the stochastic process governing these shocks to housing demand only from wage data.⁷ This is not to claim that there are no other shocks that will affect the volatility of both prices and construction. There are, but our approach still provides some quantitative measure of how misspecified our housing models would be if we were to ignore these additional shocks.

To generate predictions from the model, we need to calibrate eight parameters: $(r, \alpha, w, \delta, \theta, \sigma, c_1, c_2)$. The parameters (δ, θ, σ) govern housing demand. Consistent with the spirit of the calibration exercise described in the previous paragraph, we estimate these parameters exclusively using wage data. Identifying the remaining five parameters using only data on deviations of housing prices and construction of new houses from their steady state levels turns out to be infeasible.⁸ Therefore, we borrow estimates of the real interest rate, r , the slope of the inverse housing demand equation, α , and the slope of labor demand, w , from other sources. Finally, we use data on housing prices and quantities to estimate the parameters determining the housing supply, (c_1, c_2) .

We assume that r equals 0.04. This value is higher than standard estimates of the real interest rate because it is also meant to reflect other aspects of the cost of owning such as taxes or maintenance expenses that roughly scale up with the cost of the house. Different values of the real interest rate have little impact on our calibration, as long as it is assumed to be constant.

The value of α reflects the impact that an increase in the housing stock will have on the willingness to pay to live in a locale. If population was fixed, equation (2) would imply that the derivative of steady state housing prices with respect to the number of homes equals $-(1+r)\alpha/r$, which can be interpreted as the slope of the housing demand curve. Typically, housing demand relationships are estimated as elasticities, so we must

⁷There can still be long run trends in amenities that differ across metropolitan areas, but these will not impact the short term housing price and construction dynamics that are the focus of our simulations.

⁸As will be seen in the next Section, in order to identify the parameters of the model, we derive moment conditions from the equation in Lemma 2. More moment conditions than parameters we have to identify are derived. Nevertheless, when we try to simultaneously identify the five parameters (r, α, w, c_1, c_2) , the resulting objective function is relatively flat and identification is very weak.

first convert elasticities into the comparable slope in levels and then multiply by $r/(1+r)$. Many housing demand elasticity estimates are around one (or slightly below, in absolute value; see, e.g. Polinsky and Ellwood (1979) or Saiz (2003)), and there is a wide range in the literature, so we experiment with a range from 0 to 2. To transform the elasticity into slope in levels, we multiply by an average ratio of price to population, and that produces a range of estimates for $(1+r)\alpha/r$ ranging from 0 to 3. Multiplying this range by $r/(1+r)$ yields a range from 0 to 0.15. We use a parameter value of 0.1 in our estimation, which implies that for every 10,000 extra homes sold the marginal purchaser likes living in the area \$1,000 less per year.

Lower values do not significantly change our estimates. Even with $\alpha = 0.1$, most of the variation in house prices comes from direct shocks to wages and not from variation in congestion effects. Lemma 2 shows that we can decompose the variation in house prices from trend as

$$H_t - \hat{H}_t = \underbrace{x_t + \frac{E_t(x_{t+1})}{\bar{\phi} - \delta}}_{\text{wage shocks}} - \underbrace{\frac{\alpha(1+r)}{1+r-\phi}(N_t - \hat{N}_t)}_{\text{congestion effects}}. \quad (1.6)$$

Table 1.1 lists the volatility of each term using the parameters we estimate for each of the three regions of the United States (calculation details are in the technical appendix). In all three cases, wage shocks are much more important than variation in congestion effects. The value of α is much more important in determining the steady-state (i.e. trend) size of the city, but this steady-state is not our focus here.

The parameter α combines the impact that extra population has on wage levels with the impact that extra population has on amenities, and we also must use a distinct estimate of the connection between population and wage levels to correct our wage series for the change in population. Given the absence of compelling evidence on the links between population size and amenity levels, and the possibility that the link is actually positive (if access to other people is a consumption amenity), we make the simplifying assumption that the impact of population on amenities is zero, so that the value of α is the same as the value of α_W . While we do not literally believe this, assuming it has little impact on our estimates since it only

serves to allow us to infer productivity changes from wage changes by correcting for the changes in population. As year-to-year population changes are relatively modest, different means of correcting for population changes have little impact on the inferred productivity series.

In principle all eight parameters in our model could differ across each metropolitan area, but data limitations make it impossible for us to precisely estimate distinct values for each location. Instead, we assume the calibrated parameters (r, α, w) to be identical for all metropolitan areas and we estimate different values of the parameters $(\delta, \theta, \sigma, c_1, c_2)$ for three different regions of the U.S.⁹ Our three regions are coastal, sunbelt and interior. Metropolitan areas whose centroids are within 50 miles of the Atlantic or Pacific Oceans are defined as coastal. Metropolitan areas more than 50 miles from either coast and which are in the broad swath of southern and western states on the southern border of the country running from Florida through Arizona are defined to be in the sunbelt region. The remainder of our metropolitan areas are defined as being in the interior region of the country.

⁹Obtaining different estimates of (r, α, w) for each of these three areas is impossible, as the sources from which we borrow those estimates do not provide such detail.

Table 1.1: *Relative Volatility of Terms in House Price Equation*

	Coastal	Sunbelt	Interior
Wage Shocks	44,000	12,000	13,000
Congestion Effects	4,000	5,000	7,000

Notes: The house price equation is decomposed in equation (1.6). The volatilities are computed using the estimates in Table 1.2. Details on the computation are provided in the technical appendix.

1.3.1 Data

For our estimation exercise, we need data on housing prices, construction of new houses, number of households potentially supplying labor, and income per household for a significant number of metropolitan areas.

The housing price data is based on Federal Housing Finance Agency repeat sales indices. Construction data are housing permits reported by the U.S. Census. To estimate annual changes in the number of households, we impute the housing stock based on decadal census estimates of the housing stock and annual permits data. Specifically, we estimate the housing stock at time $t + j$ to be

$$N_t^i + \frac{\sum_{k=0}^{j-1} \text{Permits}_{t+k}^i}{\sum_{k=0}^9 \text{Permits}_{t+k}^i} (N_{t+10}^i - N_t^i),$$

where N_t^i and N_{t+10}^i are the housing stocks measured during the two closest censuses in metropolitan area i . Thus, the change in housing stock is partitioned across years based on the observed permitting activity.

Our primary source of income data comes from the Home Mortgage Disclosure Act (HMDA) files on reported income on mortgage applications. We observe all loan applicants, not just successful buyers. The HMDA data extend back to 1990. Since HMDA is essentially a 100 percent sample of everyone who sought a mortgage, the sample sizes are quite large and we have data for every metropolitan area. Importantly, the HMDA data captures household level income, which is the appropriate level given our model. The disadvantages of using HMDA income data are a relatively short time series, the fact that we do not observe those who searched but did not apply for a mortgage, and that the homebuying decision is endogenous, which can create biases because the selected sample of people who decide to apply for a loan can differ across markets or years.

An alternative data source on income is the Bureau of Economic Analysis (BEA) per capita income measure. It is available beginning in 1980 and for all metropolitan areas. However, it suffers from a number of drawbacks. First, it is at the individual, not household, level as its name implies. Households, not individuals, purchase housing units. Hence, in

our experimentation with this measure, we translate per capita incomes into household-levels by multiplying by 2.63, which is the average number of people per housing unit in our sample of areas in 1990. It also captures the incomes of many people who were not potential buyers. The incomes earned by permanent renters or people who have been immobile homeowners for many years may not have much to do with the advantage that a location brings to the marginal purchaser. In addition, the incomes of renters are both lower and less volatile than those of owners. Hence, the BEA series is likely to understate the relevant volatility in local incomes, which is critical given our purposes.¹⁰

While we experimented with both income measures, we believe the advantages of the HMDA series far outweighs its negatives. Hence, we report results using this series and comment on findings with the BEA data where appropriate.

The sample used in the estimation has 21 sunbelt metropolitan areas, 32 coastal metropolitan areas, and 60 interior ones. The data for housing prices, construction, number of households, and borrower income spans the period 1990-2004.

1.3.2 Methodology

As indicated above, we estimate the parameters $(\delta, \theta, \sigma, c_1, c_2)$ subject to particular values of (r, α, w) . We estimate these five parameters using a *sequential* two-step Generalized Method of Moments estimator.¹¹ Our two stage procedure estimates our parameters by first using

¹⁰Based on data from the New York City Housing and Vacancy Surveys (NYCHVS) from 1978-2002, the income of recent homebuyers increases by \$1.29 for every dollar increase in BEA-reported per capita income, while that for renters only rises by \$0.47. The NYCHVS only covers one city, but it highlights that the volatility of BEA per capita income is lowered by its incorporation of renter income.

¹¹The details of this estimation method are provided in the Appendix. Hansen (1982) proves consistency and asymptotic normality for the standard two-step GMM estimator, in which all parameters are simultaneously estimated. Newey (1984) expands these results and provides the correct formula for the asymptotic variance of the two-step GMM estimator of a subvector of parameters, when the moments are a function of previous GMM estimates of a different subvector of parameters. Finally, Newey and McFadden (1994) show that the sequential GMM estimators belong to the more general family of *extremum* estimators. These results guarantee that the sequential two-step GMM estimator we use is consistent, asymptotically normal and has the asymptotic variances described in the Appendix. In principle, we could estimate all of our parameters simultaneously, using information on wages, construction levels and housing prices, but, as indicated above, this would contradict the spirit of the exercise we want to perform. If we were to use data on deviations of housing prices and construction levels with respect to their steady state in order to identify the parameters (δ, θ, σ) , then our estimates of the stochastic process governing housing demand would capture not only the income process

the population-corrected wage series to estimate the housing demand parameters and then using housing price and construction series to identify the housing supply parameters. More specifically, the parameters (δ, θ, σ) are estimated from an equilibrium equation in the labor market using a two-step GMM estimator. Given these estimates, the parameters (c_1, c_2) are estimated from the equilibrium equations for the housing market in Lemma 1 using again a two-step GMM estimator.

Description of Moments

The vector of moments used to estimate (δ, θ, σ) is based on the reduced form relationship between productivity per worker and the equilibrium number of workers: $W_t^i = \widehat{W}_t^i - \alpha_W N_t^i$. The assumption that x_t^i works entirely through the wage process allows us to write: $W_t^i = w_0^i + w_1^i t + x_t^i - \alpha_W N_t^i$, which allows for a city-specific constant and a region-specific time trend in labor demand.¹² Using this expression for wages as well as the assumed value of w , we define our productivity variable, which is wages normalized for changes in the number of workers: $\widetilde{W}_t^i = W_t^i + \alpha_W N_t^i$. The resulting equation is: $\widetilde{W}_t^i = w_0^i + w_1^i t + x_t^i$, where x_t^i follows an ARMA(1,1) process. The stochastic process for the shocks is therefore

$$x_t^i = \delta x_{t-1}^i + \epsilon_t^i + \theta \epsilon_{t-1}^i,$$

with ϵ_t^i independently and identically distributed over time with

$$\mathbb{E}[\epsilon_t^i | x_t^i, x_{t-1}^i] = 0, \quad \text{and} \quad \text{var}[\epsilon_t^i | x_t^i, x_{t-1}^i] = \sigma_\epsilon^2.$$

Using these two restrictions on ϵ and data on \widetilde{W}_t^i , we identify the parameter vector (δ, θ, σ) through a vector of moments

$$\mathbb{E}[f(\widetilde{W}^i; (\delta, \theta, \sigma))] = 0.$$

(as the model indicates should be the case) but also the stochastic process governing any other unobservable variable or shock that might affect the equilibrium in the housing market.

¹²We have tried to allow for city-specific time trends but, given the short length of the time series available for estimation, this impedes the identification of the remaining parameters of the wage equation.

The exact functional form of the moment function $f(\tilde{W}^i; (\delta, \theta, \sigma))$ is contained in the Appendix. This moment function is based on different moments of the one- period changes in our productivity measure, $\Delta \tilde{W}$, and relies on the ϵ shocks having mean zero, being uncorrelated with lagged values of \tilde{W}^i , and having constant variance.¹³

Given the first stage estimates of the housing demand parameters, $(\hat{\delta}, \hat{\theta}, \hat{\sigma}^2)$, we use the equilibrium equations in Lemma 1 to build moment conditions that allow us to identify the vector (c_1, c_2) . Identification of these two parameters is performed through the vector of moment conditions:

$$\mathbb{E}[v(H^i, N^i, I^i; (c_1, c_2))] = 0.$$

The exact functional form of the moment function $v(H^i, N^i, I^i; (c_1, c_2))$ also is reported in the Appendix. This moment function is based on different moments of the deviations between the vector of housing prices, construction, and number of households and their steady state levels, $(H - \hat{H}, I - \hat{I}, N - \hat{N})$. The moments defined by the moment function $v(H^i, N^i, I^i; (c_1, c_2))$ rely on the ϵ shocks having mean zero, being uncorrelated with lagged values of N^i , and having constant variance.

In order to build the sample analogues of

$$\mathbb{E}[f(\tilde{W}^i; (\delta, \theta, \sigma))] = 0,$$

$$\mathbb{E}[v(H^i, N^i, I^i; (c_1, c_2))] = 0,$$

we use sample moment conditions that pool all the observations across metropolitan areas and time periods which we assume share the same values of the parameter vector $(\delta, \theta, \sigma, c_1, c_2)$. Specifically, we build the sample analogue of the moment conditions aggregating across metropolitan areas within regions and over our entire sample period. We pool observations across metropolitan areas, instead of splitting them across different moment conditions, to

¹³As a robustness check, we have also estimated (δ, θ, σ) using a multiple-step estimation procedure. In the first step, we use the Arellano-Bond estimator to obtain estimates of delta. Given this estimate of δ , we use a Classical Minimum Distance estimator for θ based on the first and second order temporal autocorrelation. Finally, using our estimates of (δ, θ) , we estimate σ from the residual variance. The results are very similar to the ones based on the simultaneous estimation of (δ, θ, σ) using the moment function $f(\tilde{W}^i; (\delta, \theta, \sigma))$ and are available upon request.

increase our sample size. After all, GMM estimators have optimal statistical properties only when the number of observations used in each moment condition goes to infinity, and the standard errors of our GMM estimates are valid only asymptotically.

1.3.3 Estimation Results

Table 1.2 reports our estimated parameters. The estimates of the labor demand shocks persistence parameter, δ , are 0.88 in the interior and coastal areas and 0.89 in the sunbelt. While the similarity of these estimates is striking, they are still somewhat imprecise. We cannot reject the possibility that income shocks follow a random walk (i.e., the persistence parameter equals one) and we also cannot reject much more significant mean reversion.

Table 1.2: *Estimated Demand and Supply Parameters*
HMDA Income Data, 1990- 2004

	Coastal	Sunbelt	Interior
Demand			
δ	0.88 (0.11)	0.89 (0.13)	0.88 (0.10)
θ	0.82 (0.62)	0.13 (0.13)	0.20 (0.10)
σ_ϵ	\$1,700 (500)	\$1,300 (200)	\$1,300 (100)
Supply			
c_1	10.62 (0.58)	1.47 (0.14)	3.16 (0.25)
c_2	4.08 (0.77)	0.34 (0.08)	0.12 (0.11)

Notes: δ , θ , and σ_ϵ are the autocorrelation parameter, moving average parameter and residual variance of an ARMA(1,1) estimated for the component of wages that is not explained by a linear time trend and a metropolitan area-specific constant. c_1 denotes the derivative of expected future housing prices with respect to current investment in housing construction; and c_2 denote the derivative of the physical capital cost of building a home with respect to the stock of houses. The standard errors for the demand parameters are efficient two-step GMM standard errors. The ones for the supply parameters account for error coming from the demand estimates.

The estimates of the moving average parameter θ are statistically indistinguishable from zero in the sunbelt and coastal regions. In the interior region, this moving average component estimate is 0.2 and is marginally significantly different from zero. The productivity shock estimates range from \$1,300 in the sunbelt and interior to \$1,700 on the coast. Our estimates of the housing supply parameters reported in the bottom panel of Table 1.2 indicate a value for c_1 of 10.62 in the coastal region. This implies that a 1,000 unit increase in the number of building permits in a given year raises the cost of supplying a home by \$10,620. We estimate a value of c_2 in that region of 4.08, meaning that as the number of units in a metropolitan area increases by 10,000 the cost of supplying a home increases by more than \$40,000. The estimates of c_1 are much lower in the sunbelt and interior regions, at 1.47 and 3.16, respectively. In these two regions, the estimates of c_2 are 0.34 and 0.12, respectively. Housing supply does appear to be far more elastic in those regions.¹⁴

These latter findings can be compared with the housing supply estimates reported by Topel and Rosen (1988), who use aggregate national data to estimate an elasticity of housing supply with respect to price that is between 1 and 3. In our model, that supply elasticity equals $H_t / (c_1 I_t)$. In 1990, average prices were about \$130,000. Average construction levels in a metropolitan area is approximately 8,350 units, as measured by building permits issued. If we take the Topel and Rosen (1988) elasticity to be 3, then this implies a value of c_1 of 5, which lies in the middle of our estimates.

¹⁴As noted above, we generated separate estimates using BEA per capita income data in lieu of HMDA data. This has the advantage of including years back to 1980, but we also suspect it might grossly underestimate income volatility, which is critical for our purposes. In fact, estimates of the productivity shocks are much lower, with the largest estimate of \$1,200 for coastal region markets being smaller than that reported above for sunbelt and interior markets using HMDA data. The moving average parameters are somewhat smaller across all regions, but they are also imprecisely estimated, as was the case with the estimates based on HMDA. The BEA data imply greater differences across regions in the demand shock persistence parameter, δ , with estimates ranging from 0.73 in the interior (and we can reject that coefficient equals one at standard confidence levels) to 0.8 in coastal areas and 0.9 in the sunbelt region. Estimates of supply parameters using BEA per capita income show a very similar pattern to those reported above, albeit with small point estimates. The coastal c_1 is 6.1 and its c_2 is 1.9; those for the interior and sunbelt regions are much closer to zero. See the appendix for the analogue to Table 1.4 based on using BEA per capita income in lieu of HMDA-based income.

1.4 Matching the Data and Discussion

The model presented in Section 2 implies a particular stochastic process for housing prices and for the construction of new houses. If shocks are known as they occur, then it is straightforward to show that our model implies the following ARMA(2,3) process for housing prices, with the parameter vector restricted as outlined in the appendix:

$$\Delta H_t^i = a_0^i + a_1 \Delta H_{t-1}^i + a_2 \Delta H_{t-2}^i + b_0 \epsilon_t^i + b_1 \epsilon_{t-1}^i + b_2 \epsilon_{t-2}^i + b_3 \epsilon_{t-3}^i.$$

Analogously, the model implies the following ARMA(2,1) process for the construction of new homes, with the parameter vector restricted as shown in the appendix:

$$I_t^i = d_0^i + d_1 I_{t-1}^i + d_2 I_{t-2}^i + e_0 \epsilon_{t-1}^i + e_1 \epsilon_{t-2}^i.$$

We then use these two ARMA processes, together with the estimated values of the supply and demand parameters, to derive various predictions of the model over different time horizons. Certain moments directly estimated from the data are compared to those analytically derived. In doing so, we focus on a particular set of moments of these stochastic processes: serial correlations and variances at the one, three and five year horizons. We do not focus on any contemporaneous or lagged correlations between prices and quantities for the reasons discussed next, even though much research in urban and real estate economics uses results from regressions of high frequency prices (or price changes) on demand factors such as income (or income changes).

1.4.1 The Impact of Information on the Predictions of the Model

The model discussed above assumes that shocks are observed as they occur, but we are far from confident that they are not known ahead of time. And, the results of contemporaneous correlations are sensitive to what one assumes about the underlying information structure (i.e., whether information about the change in income becomes known ahead of time or only contemporaneously with its public release). In contrast, *autocorrelations* of price and construction series are much less sensitive to information timing as we now demonstrate by

comparing the predictions of the model with our assumed information structure and the predictions if shocks are known one period ahead of time.

For this exercise, we use parameter estimates from the coastal region: $r = 0.04$, $\alpha = 0.1$, $c_1 = 10.62$, $c_2 = 4.08$, $\theta = 0.82$, $\delta = 0.88$, and $\sigma = \$1,700$. The first column in Table 1.3 reports our model's predictions for a number of variables presuming such contemporaneous knowledge.¹⁵ The second column represents our model's predictions when individuals learn about the income shock one period before it actually impacts wages.

Advance knowledge slightly increases construction volatility and adds some momentum to house price changes. Otherwise the autocorrelations are essentially unchanged. Therefore, the predictions of our model for these moments are robust to a possible misspecification of the information structure and a potential lag between the time the income shocks are known to the agents and when they are made public.

In stark contrast, the impact of the information structure on the contemporaneous correlation between changes in prices and changes in income is enormous. The bottom panel of Table 1.3 shows that if knowledge is contemporaneous to the shock, then the correlation of price and income changes over short horizons is 0.80. If individuals acquire knowledge one year ahead, then the predicted correlation is only 0.08. The correlation is only somewhat more stable at lower frequencies.

Because these correlations are so sensitive to small changes in the underlying information conditions, we focus our analysis on the serial correlation properties and volatility of price changes and construction activity.¹⁶

¹⁵For any j year interval, these predictions reflect the relationship between what happened between time t and $t - j$ and what happened between time t and $t + j$.

¹⁶Over longer horizons, a one-year shift in when information becomes known is less important, so it certainly can make good sense to explore various longer-run relationships with price changes. Because our interest is in higher frequency changes, we do not do that here.

1.4.2 Volatility and Serial Correlation in House Prices

Table 1.4 documents how well the model matches the data by comparing the model's predictions of short- and long-run volatility and serial correlation in house price changes

Table 1.3: *Sensitivity of Predictions to Different Information Structures*

Horizon	Contemporaneous Knowledge	Knowledge One Year Ahead
<i>Serial Correlation of Construction</i>		
1 year	0.51	0.56
3 year	0.18	0.19
5 year	-0.04	-0.03
<i>Volatility of Construction (units)</i>		
1 year	1,800	2,000
3 year	4,300	4,800
5 year	6,000	6,700
<i>Serial Correlation of House Price Changes</i>		
1 year	-0.00	0.09
3 year	-0.16	-0.10
5 year	-0.24	-0.21
<i>Volatility of House Price Changes (\$)</i>		
1 year	18,000	17,000
3 year	30,000	31,000
5 year	37,000	39,000
<i>Correlation of Income Changes and House Price Changes</i>		
1 year	0.80	0.08
3 year	0.93	0.61
5 year	0.95	0.75

Notes: The parameter values estimated for the coastal region using HMDA wage data are assumed here: $\delta = 0.88$, $\theta = 0.82$, $\sigma_\epsilon = \$1,700$, $c_1 = 10.62$, and $c_2 = 4.08$.

and new construction with the actual moments from the data. Standard deviations and serial correlation coefficients from the underlying data over this time period are reported in columns adjacent to our model predictions.

Table 1.4: *Volatility and Serial Correlation in House Prices and Construction: HMDA Income Data, 1990-2004*

Horizon	Coastal		Sunbelt		Interior	
	Model	Data	Model	Data	Model	Data
<i>Volatility of House Price Changes (\$)</i>						
1 year	18,000	13,300	5,000	2,000	6,000	3,600
3 year	30,000	34,100	8,000	4,400	10,000	8,400
5 year	37,000	48,300	9,000	5,400	12,000	10,700
<i>Serial Correlation of House Price Changes</i>						
1 year	-0.00	0.84	-0.12	0.64	-0.07	0.73
3 year	-0.16	0.32	-0.28	-0.09	-0.25	0.10
5 year	-0.24	-0.80	-0.35	-0.73	-0.36	-0.72
<i>Volatility of Construction (units)</i>						
1 year	1,800	1,900	3,600	5,300	2,000	1,600
3 year	4,200	4,600	9,000	12,000	5,700	3,800
5 year	5,900	6,300	12,000	15,500	8,600	5,000
<i>Serial Correlation of Construction</i>						
1 year	0.50	0.75	0.56	0.82	0.72	0.74
3 year	0.17	0.18	0.25	0.23	0.47	0.25
5 year	-0.04	-0.79	0.03	-0.60	0.25	-0.72

Notes: The moments computed from the data allows the mean of housing price changes and construction to vary across metropolitan areas. The moments generated from the model use the estimates in Table 1.2.

Volatility in House Prices

The model generally overpredicts price volatility except in the coastal region at 3- and 5-year horizons. One explanation for this excess predicted volatility is that the HMDA data may be overestimating the actual volatility in local labor demand. Predicted volatility is closer to the data in both absolute and percentage terms over longer horizons in the interior regions. Those differences are within \$2,000. And, the model captures the sharply rising volatility in price changes over longer horizons in coastal markets,¹⁷ but it never matches the very high price volatility seen in those areas over 3- and 5-year horizons. Except in coastal markets, there appears to be more than enough volatility in local income processes to account for house price volatility.¹⁸

Serial Correlation in House Prices

Turning now to the model predictions about the serial correlation of house price changes over 1, 3 and 5 year horizons reported in the second panel of Table 1.4, the model predicts very modest autocorrelation of one- year price changes, ranging from zero in the coastal region to -0.12 for the sunbelt region. Comparing these predictions with the actual data reveals a glaring mismatch between the model and reality. In the real world, as Case and Shiller (1989) documented long ago, there is strong positive serial correlation at one-year frequencies. A one dollar increase in prices during one year is associated with between a 64 and 84 cent increase in prices during the next period, depending upon region.

There is no reasonable calibration of the model that can match the strong positive serial correlation of prices at high frequencies. One possible explanation lies in the microfoundations of the housing market. If there is a learning process at work, whereby people

¹⁷This is due to the higher underlying volatility in the local income process (σ is 30 percent higher in the coastal metropolitan areas), as well as higher moving average component θ .

¹⁸The results are far different if the BEA income series is used. In that case, the model grossly underpredicts price change variation, by 50%-75% or more. See the appendix for the analogue to Table 1.4 based on BEA per capita income. Thus, if one disagrees with our conclusion that the HMDA-based income series is superior and that per capita income better reflects reality, then local housing markets are far too volatile given their (income) fundamentals.

gradually infer the state of demand from prices, then this can generate serial correlation. An alternative explanation is less rational: people see past price changes and infer future price growth (as in Glaeser *et al.* (2008)). Neither idea is captured in our model. In our model, individuals are fully rational and they know the parameters that govern the stochastic process for housing prices and construction of new houses.

At three year periods, the model and the data continue to diverge. The model continues to predict mean reversion in prices, with the implied serial correlation coefficient ranging from -0.16 for the coastal region to -0.28 for the sunbelt region. The real data shows at least mild positive serial correlation for all but the sunbelt region. Once again, price changes are too positively correlated to match the model.

At 5-year time horizons, the model correctly predicts that price changes mean revert, which is an important stylized fact about local housing markets. However, the point estimates are well below the amount of mean reversion apparent in the data. This is one case in which we are skeptical of the data because our procedures for detrending, which involve subtracting the metro area means, probably induce some spurious mean reversion given the limited fifteen year time series.

While part of the reason for the magnitude mismatch may be due to this factor, that does not provide a complete explanation. If we lengthen the price change time series and include the 1980s, computed mean reversion is lower, but is still higher than our estimates in Table 1.4 . For example, the serial correlation in five year price changes falls from - 0.80 to -0.57 in the coastal region. That still is more than double the -0.24 estimate yielded by our model (Table 1.4). And, using BEA per capita income over the longer time period dating back to 1980 does not yield a perfect (or close to perfect) match either.¹⁹ Hence, the model should be viewed as successful in capturing the fact that there is mean reversion in price changes over long horizons, but it fails to match the strength of that pattern.

¹⁹Similar patterns are evident in the other regions.

1.4.3 Volatility and Serial Correlation in Construction

Volatility in Construction

The model matches the volatility of construction activity at all time horizons in the coastal region quite well, and especially at high frequency (panel 3, Table 1.4). The match quality is less good, but tolerable, in the sunbelt region. The model predicts much greater volatility over longer horizons, but underpredicts volatility by one- quarter to one-third in this region. We consistently overpredict construction quantity by at least 25% at each horizon in interior markets.²⁰

Serial Correlation in Construction

In stark contrast to the model's complete failure to predict strong persistence in price changes over one-year horizons, it always correctly predicts positive, high frequency serial correlation in construction in all regions, with the match being very good for the interior region. Our estimates are about one-third below what the actual data show for the coastal and sunbelt regions, so complete success for the model cannot be claimed here. We do better at 3-year horizons. Our model estimates correctly mimic the lower level of serial correlation at this longer horizon in all regions. And, our point estimates are very close matches to the data in the coastal and sunbelt regions.

However, the estimates over 5-year horizons do not match the data. As noted above, we are skeptical of the value of creating such differences using only 15 years of data. If we go back and include the 1980s, calculated mean reversion fall by about two-thirds in each region (e.g., from -0.79 to -0.27 in the coastal region; from -0.60 to -0.20 in the sunbelt region; and from -0.72 to -0.24 in the interior region). Thus, it certainly looks as if the short time span over which we have higher quality income data is leading to an upwardly biased level of mean reversion in construction for the model to match. That said, our model estimates

²⁰As was the case for price change volatility, using per capita income from the BEA in lieu of household-level income from HMDA leads us to dramatically underpredict construction volatility. To reiterate, if one believes the BEA series more accurately reflects the true variation of local income processes, then housing markets are far too volatile relative to their fundamentals.

still do not match those lower levels of mean reversion.²¹

1.5 Conclusion

This paper presents a dynamic linear rational expectations model of housing markets based on cross-city spatial equilibrium conditions. Its aim is to show how well a housing model that focuses on income shocks may approximate certain features of the housing market. The model predicts that housing markets will be largely local, which they are, and that construction persistence is fully compatible with price mean-reversion. The model is also consistent with price changes being predictable.

The model has notable successes and failures at fitting the real data. It generally captures important differences across types of markets, especially coastal ones that have inelastic supply sides to their housing markets. The model also does a decent job of accounting for variation in price changes. An important implicit assumption underlying that conclusion is that the HMDA series more accurately reflects the volatility of local income processes than (say) the BEA's per capita income measure. More in-depth research on this data issue seems warranted given its importance in allowing the model to approximate market price volatility. This conclusion also generally applies to the volatility of quantities as reflected in construction permits.

That said, we still cannot precisely match the very high volatility of three- and five-year price changes observed in the inelastically supplied coastal regions. Thus, it also would be useful for future research to try to pin down whether there is excess volatility in those markets.

The model does tolerably well at accounting for the strong positive serial correlation of construction quantities from one year to the next. It also correctly captures the weakening of this persistence over longer horizons, but fails to match the magnitude of the mean reversion

²¹This is the one case in which using the BEA data on income and the longer time series including the 1980s leads to better matches. In this case, the model always predicts at least modest mean reversion in construction over 5-year horizons, and the match quality is quite good for the interior region.

in quantities over longer horizons especially. Some of the failure in matching the magnitude of mean reversion in prices and quantities over longer horizons may be due to data error, but that is not a complete explanation. This is another avenue for fruitful research.

The model fails utterly at explaining the strong, high frequency positive serial correlation of price changes. It does a much better job of accounting for the mean reversion over longer, five-year horizons, especially when one takes into account the likelihood our procedures overstate true mean reversion over this longer time span.

This suggests that housing economists have one very big puzzle to explain, along with some other issues. The major puzzle is the strong persistence in high frequency price changes from one year to the next. This failure must be viewed as stark given that attempt to match moments for a time period that does not include the recent extraordinary boom and bust. Other matters that certainly merit closer scrutiny include the extremely high price change volatility in coast markets over longer time horizons and the inability to match mean reversion in construction over longer horizons. These empirical misses are significant, but it remains true that a dynamic urban model can account for many of the important features of housing markets. We see this model as a starting point for a larger agenda of research on real estate dynamics that starts with a dynamic spatial equilibrium model. One natural extension is to include interest rate volatility, and we have sketched such an approach in an earlier version of this paper. A second extension is to relax the assumption of perfect rationality for home-buyers, and perhaps builders as well.

Chapter 2

Arrested Development: Theory and Evidence of Supply-Side Speculation in the Housing Market¹

2.1 Introduction

How do prices aggregate information? We take up this question in a setting of particular macroeconomic importance: housing markets. Housing is a key driver of the business cycle (Leamer, 2007), and the causes of the financial crisis of 2008 and the Great Recession originated in housing markets (Mian and Sufi, 2009, 2011). An enduring feature of these markets is booms and busts in prices that coincide with widespread disagreement about fundamentals (Shiller, 2005). This paper argues that these cycles are caused by how housing markets aggregate beliefs.

Studying belief aggregation allows us to address some of the most puzzling aspects of the U.S. housing boom that occurred between 2000 and 2006. According to the standard model of housing markets, elastic housing supply prevents house price booms by allowing

¹This chapter is co-authored with Eric Zwick.

new construction to absorb rising demand.² But the episode from 2000 to 2006 witnessed several major anomalies, in which historically elastic cities experienced house price booms despite continuing to build housing rapidly. And house prices rose more in many of these cities—located in Arizona, Nevada, inland California, and Florida—than in cities where it was difficult to build new housing. Further complicating the puzzle, house prices remained flat in other elastic cities that were also rapidly building housing. Why was rapid construction able to hold down house prices in some cities and not others?

We solve this puzzle by adding two ingredients to the standard model. The first is a friction that makes owner-occupancy more efficient than renting. The second is disagreement about long-run growth paths. In this framework the way housing markets aggregate beliefs depends on a city's land availability. Prices appear more optimistic when land is plentiful and building houses is easy, reversing the standard model's intuition for how land supply influences prices. Crucially, optimism amplifies prices most when a city nears but has not yet reached a long-run development constraint. This mechanism matches the data. The anomalous cities are those that, as the boom began, found themselves in just this state of "arrested development."

We model a city of developers and residents with a fixed amount of land available for development. Developers decide how many houses to build and how much land to buy. Residents decide how much housing to consume and whether to buy or rent. They prefer owning their houses over renting because of frictions in the rental market.³ Residents can invest in the equity of developers, which provides exposure to land prices. Short-selling land and housing is impossible, but residents can short-sell developer equity. Over time, new residents arrive in the city, leading developers to build houses using their holdings of undeveloped land. Because of this growth, the city gradually exhausts its land supply. What today's investors believe about future inflows determines the price of undeveloped land.

²See, for example, Glaeser *et al.* (2008), Gyourko (2009), and Saiz (2010).

³Such frictions include the effort spent monitoring tenants to prevent property damage (Henderson and Ioannides, 1983), tax disadvantages (Poterba, 1984), and difficulty renting properties like single-family homes that are designed for owners (Glaeser and Gyourko, 2007).

House construction is instantaneous and developers bear a constant unit cost per house. As a result, all variations in house prices are caused by movements in land prices and not construction costs. Data from the U.S. boom support this feature of the model. Rapidly rising land prices account for most of the house price increases across cities. In contrast, construction costs remained relatively stable throughout the boom, and cost changes hardly varied across cities. These aspects of the data distinguish our theory from those that stress “time-to-build” factors such as input shortages or delivery lags (Mayer and Somerville, 2000; Gao, 2014).

We study a demand shock that raises the current inflow of new residents and also creates uncertainty about future inflows. Disagreement about long-run demand leads to disagreement about future house prices. The most optimistic residents seek to speculate through buying housing and through buying the equity of optimistic developers who are buying land.

Our first result is that speculation is crowded out of the housing market and into the land market. Consider an optimistic resident who wishes to speculate on future house prices. Buying a house and renting it out is difficult because of the widespread preference for owner-occupancy. And buying more housing for personal consumption is unappealing because of diminishing marginal utility. Land however offers a pure, frictionless bet on real estate. The optimistic resident chooses to invest in land through buying developer equity.

With data from the U.S. housing boom, we confirm several of the model’s predictions about land speculation. In the model, developers run by optimistic CEOs use resident financing to amass large land portfolios, buying land from less optimistic developers. Consistent with this prediction, we find that supply-side speculation figures prominently in the data. Between 2000 and 2006, the eight largest U.S. public homebuilders tripled their land investments, an increase far exceeding their additional construction needs. Their market equity then fell 74%, with most of the losses coming from write-downs on their land portfolios. The model also predicts that short-selling of developer equity increases during a boom because pessimistic residents disagree with the high valuations of the developer land

portfolios. Matching this prediction, the short interest in homebuilder stocks rose from 2% in 2001 to 12% in 2006. Rising short interest provides evidence of disagreement over the value of homebuilder land portfolios and thus over future house prices.

Our second result concerns how house prices aggregate beliefs. Speculators are crowded into the land market, while homeownership remains dispersed among residents of all beliefs. Therefore, house prices reflect a weighted average of the optimistic belief of speculators and the average owner-occupant belief. The weight on the optimistic belief equals the share of the housing market on the margin that consists of the land market. Prices look most optimistic where land is plentiful and building easy—that is, in cities where the short-run elasticity of housing supply is large.

This optimism bias affects prices most when the city's housing supply will become inelastic soon. This observation, which constitutes our third result, explains why house price booms occur in some elastic cities and not others. Consider a city in which the land available for development is large relative to the city's current size. Here, new construction fully absorbs the demand shock now and in the foreseeable future, and so beliefs about future house prices remain unchanged. The shock raises future price expectations only in cities where construction will be difficult in the near future.

Speculation amplifies house price booms most in cities that exist in a state of arrested development: they have ample land for construction today, but also face land barriers that will restrict growth in the near future. This theoretical supply condition characterizes the anomalous elastic cities during the U.S. housing boom. For instance, Las Vegas faces a development boundary put in place by Congress in 1998 and depicted in Figure 2.1. During the 2000-2006 housing boom, many investors believed the city would soon run out of land.⁴ Likewise, Phoenix's long-run development is constrained by Indian reservations and

⁴Las Vegas provides a particularly clear illustration of our model. The ample raw land available in the short-run allowed Las Vegas to build more houses per capita than any other large city in the U.S during the boom. At the same time, speculation in the land markets caused land prices to quadruple between 2000 and 2006, rising from \$150,000 per acre to \$650,000 per acre, and then lose those gains. This in turn led to a boom and bust in house prices. The high price of \$150,000 for desert land before the boom and after the bust demonstrates the binding nature of the city's long-run development constraint. A *New York Times* article published in 2007 cites investors who believed the remaining land would be fully developed by 2017 (McKinley and Palmer, 2007).

National Forests that surround the metropolitan area (Land Advisors). In inland California, much of the farmland around cities is protected by a state law that penalizes real estate development on these parcels (Onsted, 2009).

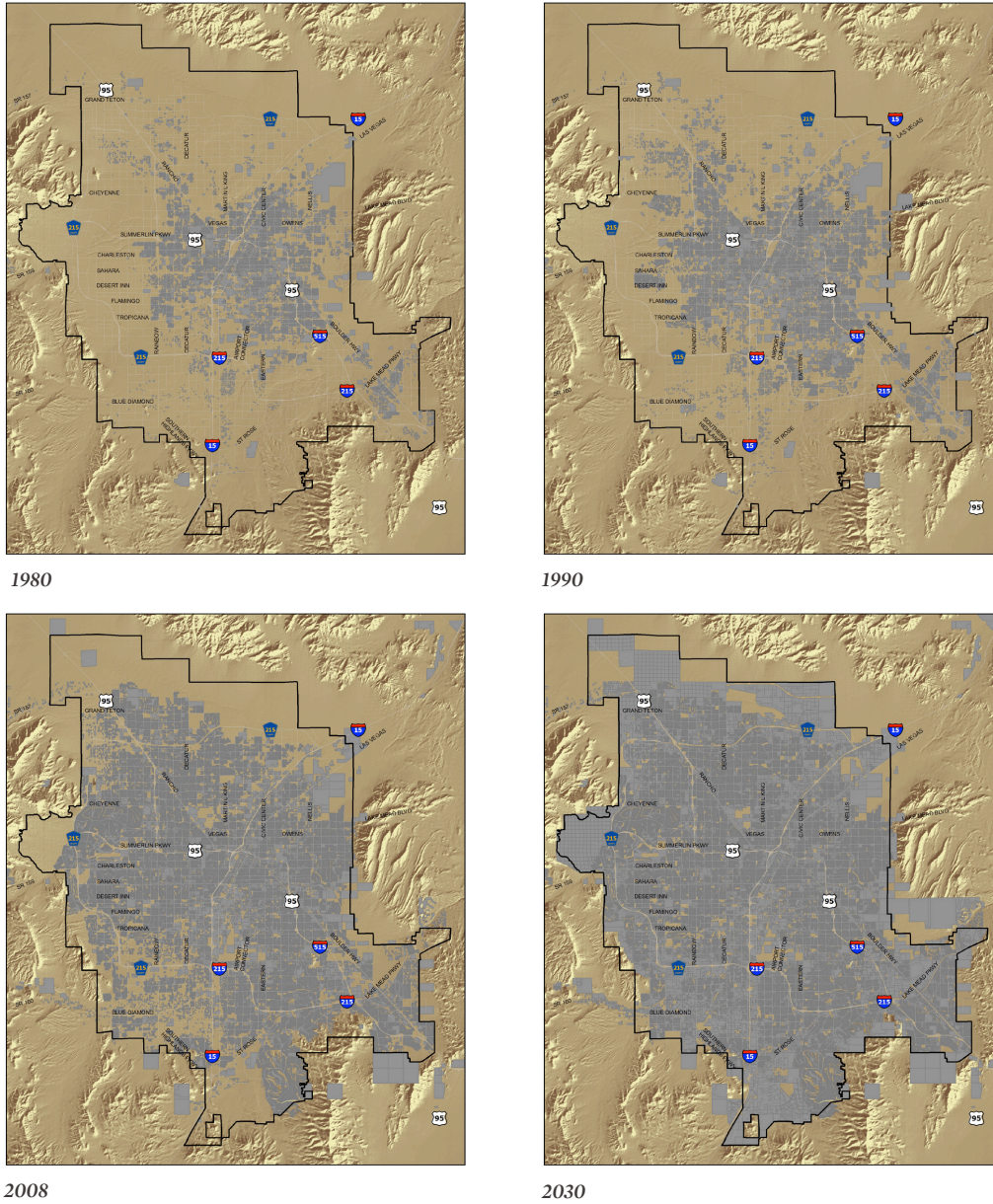
When disagreement is strong enough, house prices increase more in these nearly developed cities than in a fully developed city. In the nearly developed cities, the extreme optimistic beliefs of land speculators determine house prices, amplifying the house price boom. Prices remain more stable in the fully developed city because they reflect the average belief. This result explains the puzzling house price booms in elastic areas that motivate this paper. Supply conditions in these places—elastic current supply, inelastic long-run supply—lead disagreement to have the largest possible amplification effect on a house price boom.

Our theory differs from several other explanations for the strong house price booms that occurred in elastic areas between 2000 and 2006. One possibility is that these cities experienced much larger demand shocks than the rest of the United States.⁵ Our analysis assumes a constant demand shock across cities; the heterogeneity in city house prices booms results entirely from differences in supply conditions. An additional possibility is that uncertainty increased land values due to the embedded option to develop land with different types of housing (Titman, 1983; Grenadier, 1996), and that this option value increase was largest in cities with an intermediate amount of land. In our model, all housing is identical, so this option does not exist. A final explanation is that developers hoarded land to gain monopoly power, and the incentive to do so was strongest in cities about to run out of land. This effect does not appear in our model because homebuilding is perfectly

The dramatic rise in land prices during the boom resulted from optimistic developers taking large positions in the land market. In a striking example of supply-side speculation, a single land development fund, Focus Property Group, outbid all other firms in every large parcel land auction between 2001 and 2005 conducted by the federal government in Las Vegas, obtaining a 5% stake in the undeveloped land within the barrier. Focus Property Group declared bankruptcy in 2009.

⁵For instance, the expansion of credit described by Mian and Sufi (2009) may have been largest in these cities. Alternatively, historical increase in house prices in nearby areas may have spread to these cities, either through behavioral contagion (DeFusco *et al.*, 2013) or long-distance gentrification (Guerrieri *et al.*, 2013).

Figure 2.1: Long-Run Development Constraints in Las Vegas



Notes: This figure comes from Page 51 of the Regional Transportation Commission of Southern Nevada's Regional Transportation Plan 2009-2035 (RTCSNV). The first three pictures display the Las Vegas metropolitan area in 1980, 1990, and 2008. The final picture represents the Regional Transportation Commission's forecast for 2030. The boundary is the development barrier stipulated by the Southern Nevada Public Land Management Act. The shaded gray region denotes developed land.

competitive, as is the case empirically at the metro-area level.⁶ Unlike these stories, our approach explores the cross-sectional implications of disagreement, an under-studied aspect of housing cycles for which we provide direct evidence.

In addition to explaining the city-level cross-section, our model offers new predictions on the cross-section of neighborhoods within a city. We allow some residents to prefer renting over owner-occupancy, so that both rental and owner-occupied housing exist in equilibrium. During periods of disagreement, optimistic speculators hold the rental housing, just as they hold land. Prices appear more optimistic, and hence house price booms are larger, in neighborhoods where a greater share of housing is rented. This prediction matches the data: house prices increased more from 2000 to 2006 in neighborhoods where the share of rental housing in 2000 was higher.

A long literature in macroeconomics and finance has studied how prices aggregate information. When markets are complete and investors share a common prior, prices usually are efficient and reflect the information of all market participants (Fama, 1970; Grossman, 1976; Hellwig, 1980). Our paper sits among a body of work showing that prices reflect only a limited and potentially biased subset of information when investors persistently disagree with each other, and markets are incomplete. Many of these papers focus on strategic considerations that arise in this setting, and the implications for asset prices (Harrison and Kreps, 1978; Scheinkman and Xiong, 2003). A related literature, starting with Miller (1977), demonstrates that prices can be biased even in the absence of strategic considerations because optimists end up holding the asset.⁷ We show that this optimism bias is strongest in housing markets when land is plentiful or when much of the housing

⁶Somerville (1999) demonstrates the high level of homebuilder competition at the metro-area level, although he points out that construction is less competitive at the neighborhood level. Hoberg and Phillips (2010) argue that price booms often occur in competitive industries because firms mistakenly believe they will obtain future monopoly power.

⁷In these papers, all market participants are fundamental investors who ignore other investors' beliefs (Chen *et al.*, 2002; Geanakoplos, 2009; Hong and Sraer, 2012; Simsek, 2013a,b). Pástor and Veronesi (2003, 2009) also study environments in which investors care only about long-run fundamentals during booms and busts, but their focus is on learning, and all investors agree as they are all identical. Piazzesi and Schneider (2009) and Burnside *et al.* (2013) also apply models of disagreement to the housing market. Papers in which strategic behavior matters include Abreu and Brunnermeier (2003), Allen *et al.* (2006), and Hong *et al.* (2006).

stock is rented. In contrast, prices aggregate beliefs well in cities where the housing stock is fixed and owner-occupied. In these areas, house prices reflect the average of all resident beliefs, even though they are agreeing to disagree and short-selling housing is impossible.

The paper proceeds as follows. In Section 2.2, we document the puzzling aspects of the cross-section of the U.S. housing boom, as well as the importance of supply-side speculation in land markets. Section 2.3 models the housing market environment. Section 2.4 contains our analysis of how house prices aggregate beliefs. In Section 2.5, we derive implications of the model to explain the empirical cross-section of housing markets during the U.S. boom. Section 2.6 contains new predictions on the cross-section of neighborhoods within a city, and Section 2.7 concludes.

2.2 Stylized Facts of the U.S. Housing Boom and Bust

2.2.1 The Cross-Section of Cities

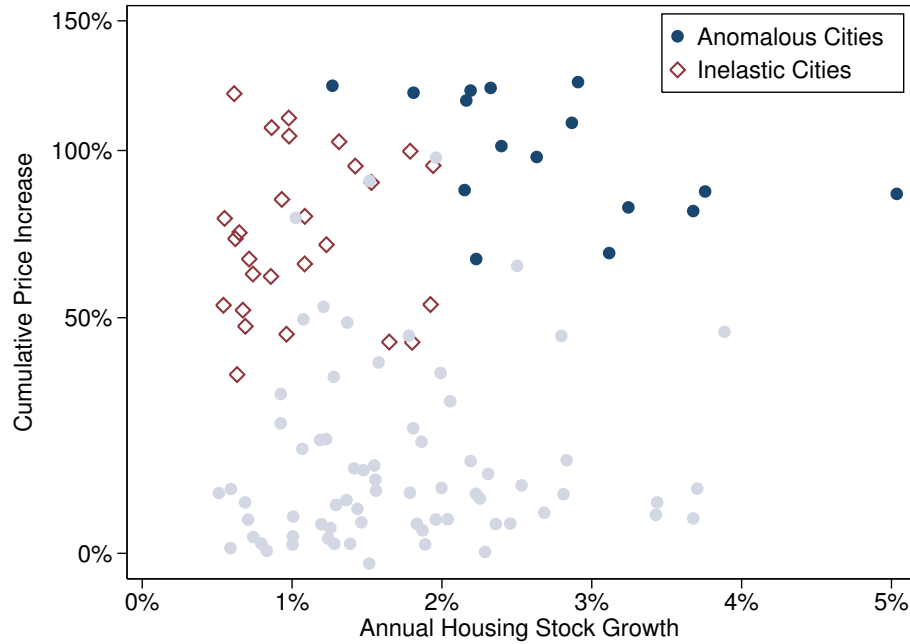
The Introduction mentions three puzzles about the cross-section of city experiences during the boom. First, large house price booms occurred in elastic cities where new construction historically had kept prices low. Second, the price booms in these elastic areas were as large, if not larger, than those happening in inelastic cities at the same time. Finally, house prices remained flat in other elastic cities that were also rapidly building housing.

We document these puzzles using city-level house price and construction data. House price data come from the Federal Housing Finance Agency's metropolitan statistical area quarterly house price indices. We measure the housing stock in each city at an annual frequency by interpolating the U.S. Census's decadal housing stock estimates with its annual housing permit figures. Throughout, we focus on the 115 metropolitan areas for which the population in 2000 exceeds 500,000. The boom consists of the period between 2000 and 2006, matching the convention in the literature to use 2006 as the end point (Mian *et al.*, 2013).

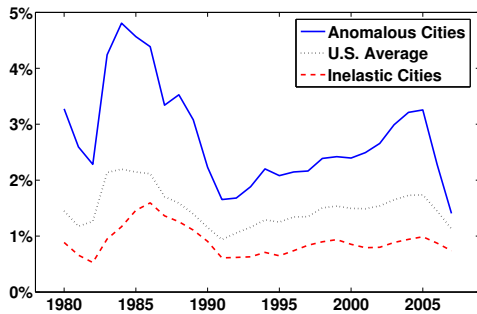
Figure 2.2(a) plots construction and house price increases across cities during the boom. The house price increases vary enormously across cities, ranging from 0% to 125% over this

Figure 2.2: The U.S. Housing Boom and Bust Across Cities

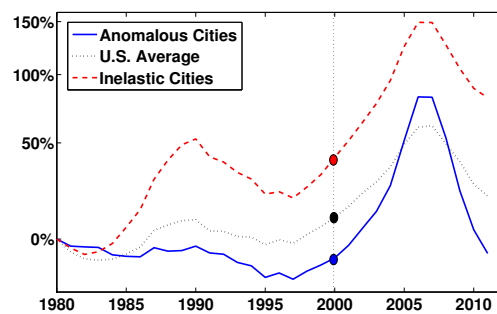
a) Price Increases and Construction, 2000-2006



b) Historic Construction



c) Historic Prices



Notes: Anomalous Cities include those in Arizona, Nevada, Florida, and inland California. Inelastic Cities are Boston, Providence, New York, Philadelphia, and all cities on the west coast of the United States. We measure the housing stock in each city at an annual frequency by interpolating the U.S. Census's decadal housing stock estimates with its annual housing permit figures. House price data come from the second quarter FHFA house price index deflated by the CPI-U. The figure includes all metropolitan areas with populations over 500,000 in 2000 for which we have data. (a) The cumulative price increase is the ratio of the house price in 2006 to the house price in 2000. The annual housing stock growth is the log difference in the housing stock in 2006 and 2000 divided by six. (b), (c) Each series is an average over cities in a group weighted by the city's housing stock in 2000. Construction is annual permitting as a fraction of the housing stock. Prices represent the cumulative returns from 1980 on the housing in each group.

brief six-year period. The largest price increases occurred in two groups of cities. The first group, which we call the Anomalous Cities, consists of Arizona, Nevada, Florida, and inland California. The other large price booms happened in the Inelastic Cities, which comprise Boston, Providence, New York, Philadelphia, and the west coast of the United States.

The history of construction and house prices in the Anomalous Cities before 2000 constitute the first puzzle. As shown in Figures 2.2(b) and 2.2(c), from 1980 to 2000 these cities provided clear examples of elastic housing markets in which prices stay low through rapid construction activity. Construction far outpaced the U.S. average while house prices remained constant. The standard model of housing cycles would have predicted the surge in U.S. housing demand between 2000 and 2006 to increase construction in these cities but not to raise prices. Empirically, the shock did increase construction, as shown in Panel (b). The puzzle is that house prices rapidly increased as well.

The second puzzle is that the price increases in the Anomalous Cities were as large as those in the Inelastic Cities. The Inelastic Cities consist of markets where house prices rise because regulation prohibits construction from absorbing higher demand. We document this relationship in Panels (b) and (c) of Figure 2.2, which show that construction in these cities was lower than the U.S. average before 2000 while house price growth greatly exceeded the U.S. average. The standard housing cycle model would have predicted the Inelastic Cities to lead the nation in house price growth in the boom after 2000. Although house prices did sharply rise, the price increases in the Inelastic Cities were no larger than those in the Anomalous Cities where the boom led to rapid construction.

The final puzzle is that some elastic cities built housing quickly during the boom but, unlike the Anomalous Cities, experienced stable house prices. These cities appear in the bottom-right corner of Figure 2.2(a), and are located in the southeastern United States (e.g. Texas and North Carolina). Their construction during the boom quantitatively matches that in the Anomalous Cities, but the price changes are significantly smaller. Why was rapid construction able to hold down house prices in some cities and not others?

One response to these three puzzles is that the Anomalous Cities simply experienced

much larger demand shocks than the rest of the nation during the boom. Although differential demand shocks surely explain part of the cross-section, they cannot account for all aspects of the Anomalous Cities just documented. These cities had been experiencing abnormally large demand shocks for years before 2000. Figure 2.2(b) shows that they were some of the fastest growing cities in the United States. Yet the surging demand to live in these areas did not increase prices. The departure from this pattern after 2000 requires a more nuanced theory than the hypothesis that housing demand increased particularly strongly in the Anomalous cities during the boom.

2.2.2 The Central Importance of Land Prices

This paper argues that speculation in land markets explains the variation in the house price boom across cities just documented. Our model demonstrates that land market speculation amplifies house price increases by making prices look more optimistic, and that this amplification is strongest in areas at the same level of development as the Anomalous Cities. In our framework, all movements in house prices arise from changes in land prices that reflect optimistic beliefs. Matching this premise, land price increases empirically account for nearly all of the increase in house prices during the boom, as we now show.

Tracing house price increases to land prices distinguishes our argument from “time-to-build” theories. According to the time-to-build hypothesis, house prices rise during a boom because of a temporary failure of homebuilders to expand construction. This delivery lag derives from obstacles erected by local regulators or from temporary shortages of inputs such as drywall and skilled labor. Under this theory, the price of undeveloped land should remain constant during the boom. Because land prices reflect the long-run, temporary housing shortages have no effect on the price of undeveloped land. These shortages instead raise construction costs and the shadow price of regulatory building permission.

To assess the importance of land prices, we gather data on land prices and construction costs at the city level. Data on land prices come from the indices developed by Nichols *et al.* (2010) using land parcel transaction data. They run hedonic regressions to control for

parcel characteristics and then derive city-level indices from the coefficients on city-specific time dummies. We measure construction costs using the R.S. Means construction cost survey. This survey asks homebuilders in each city to report the marginal cost of building a square foot of housing, including all labor and materials costs. Survey responses reflect real differences across cities in construction costs. In 2000, the lowest cost is \$54 per square foot and the highest is \$95; the mean is \$67 per square foot and the standard deviation is \$9.

Competition among homebuilders implies that, when construction is positive, house prices must equal land prices plus construction costs: $p_t^h = p_t^l + K_t$. Log-differencing this equation between 2000 and 2006 yields

$$\Delta \log p^h = \alpha \Delta \log p^l + (1 - \alpha) \Delta \log K,$$

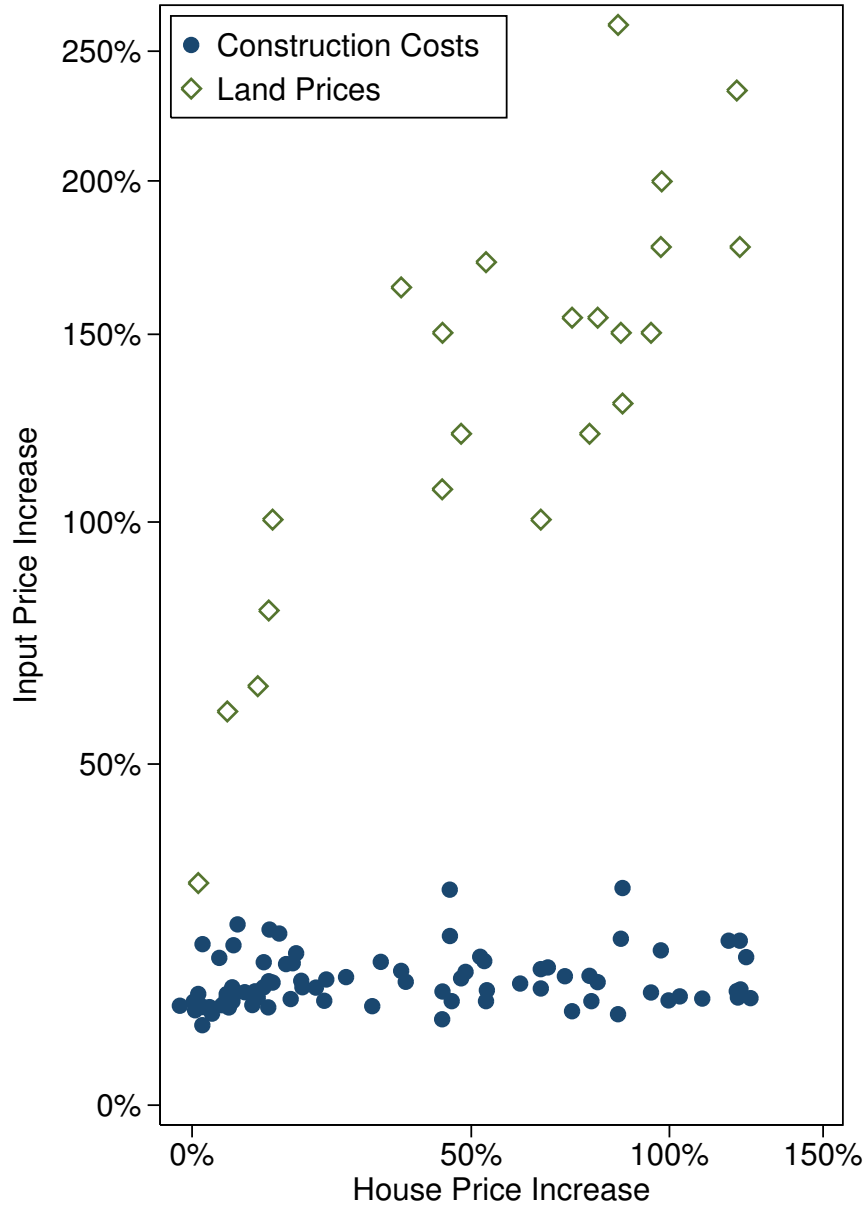
where Δ denotes the difference between 2000 and 2006 and α is land's share of house prices in 2000. The factor that matters more should vary more closely with house prices across cities. Because α and $1 - \alpha$ are less than 1, the critical factor should also rise more than house prices do.

Figure 2.3 plots for each city the real growth in construction costs and land prices between 2000 and 2006 against the corresponding growth in house prices. Construction costs rose relatively little during this period, and growth in these costs does not vary in relation to the size of house price increases. Land prices display the opposite pattern, rising substantially during the boom and exhibiting a high correlation with house prices. Each city's land price increase also exceeds its house price increase. This evidence underscores the central importance of land prices for understanding the cross-section of house price booms.

2.2.3 Land Market Speculation by Homebuilders

The land price booms just documented were driven by speculation in land markets. The term "speculation" refers to the process in which optimists buy up an asset that cannot be shorted, biasing its price. Our model describes two implications of this behavior. First,

Figure 2.3: *Input Price and House Price Increases Across Cities, 2000-2006*



Notes: We measure construction costs for each city using the R.S. Means survey figures for the marginal cost of a square foot of an average quality home, deflated by the CPI-U. Gyourko and Saiz (2006) contains further information on the survey. Land price changes come from the hedonic indices calculated in Nichols et al. (2010) using land parcel transactions, and house prices come from the second quarter FHEA housing price index deflated by the CPI-U. The figure includes all metropolitan areas with populations over 500,000 in 2000 for which we have data.

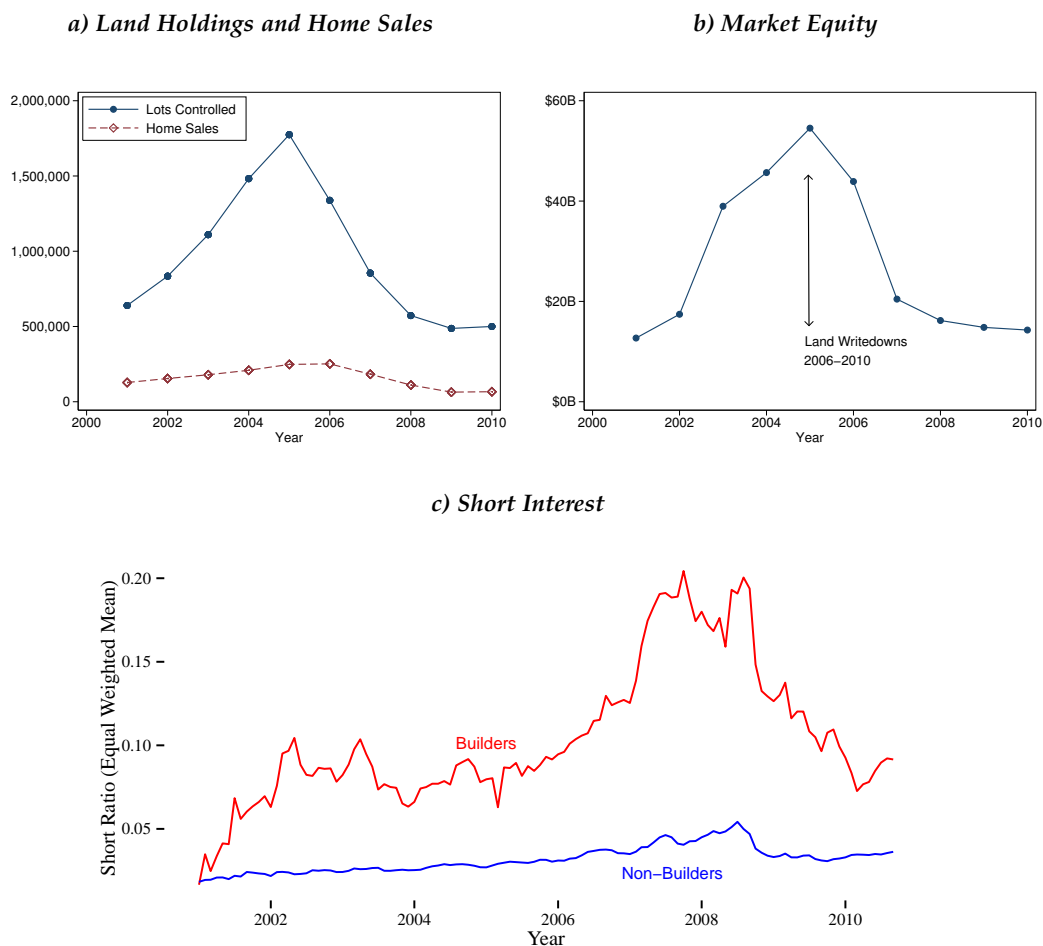
the owners of the land during the boom increase their positions as they crowd out less optimistic landowners. Second, when their beliefs are revealed to be more optimistic than reality, optimists suffer capital losses. We document both of these features among a class of landowners for whom rich data are publicly available: public homebuilders. We focus on the eight largest firms and hand-collect landholding data from their annual financial statements between 2001 and 2010.

Consistent with speculative behavior, these firms nearly tripled their landholdings between 2001 and 2005, as shown in Figure 2.4(a). These land acquisitions far exceed additional land needed for new construction. Annual home sales increased by 120,000 between 2001 and 2005, while landholdings increased by 1,100,000 lots. One lot can produce one house, so landholdings rose more than *nine times* relative to home sales. In 2005, Pulte changed the description of its business in its 10-K to say, “We consider land acquisition one of our core competencies.” This language appeared until 2008, when it was replaced by, “Homebuilding operations represent our core business.”

Having amassed large land portfolios, these firms subsequently suffered large capital losses. Figure 2.4(b) documents the dramatic rise and fall in the total market equity of these homebuilders between 2001 and 2010. Homebuilder stocks rose 430% and then fell 74% over this period. The majority of the losses borne by homebuilders arose from losses on the land portfolios they accumulated from 2001 to 2005. In 2006, these firms began reporting write-downs to their land portfolios. At \$29 billion, the value of the land losses between 2006 and 2010 accounts for 73% of the market equity losses over this time period. The homebuilders bore the entirety of their land portfolio losses. Their land acquisitions represented their beliefs about long-term land prices, as opposed to a no-downside bets made with access to free credit.

Further evidence of homebuilder optimism comes from short-selling of their market equity. If the homebuilders buying land are more optimistic than most investors, then other investors should bet against them by shorting their stock. Figure 2.4(c) plots monthly short interest ratios, defined as the ratio of shares currently sold short to total shares outstanding,

Figure 2.4: Supply-Side Speculation Among U.S. Public Homebuilders, 2001-2010



Notes: **(a), (b)** Data come from the 10-K filings of Centex, Pulte, Lennar, D.R. Horton, K.B. Homes, Toll Brothers, Hovnanian, and Southern Pacific, the eight largest public U.S. homebuilders in 2001. “Lots Controlled” equals the sum of lots directly owned and those controlled by option contracts. The cumulative writedowns to land holdings between 2006 and 2010 among these homebuilders totals \$29 billion. **(c)** Short interest is computed as the ratio of shares currently sold short to total shares outstanding. Monthly data series for shares short come from COMPUSTAT and for shares outstanding come from CRSP. Builder stocks are classified as those with NAICS code 236117.

for homebuilder stocks and non-homebuilder stocks between 2001 and 2010. Throughout the boom, short interest of homebuilder stock sharply increased, rising from 2% in 2001 to 12% in 2006. It further increased as homebuilders began to announce their land losses in 2006. Rising short interest provides direct evidence of disagreement over the value of homebuilder land portfolios and thus over future house prices.

2.3 A Housing Market with Homeowners and Developers

Housing Supply. The city we study has a fixed amount of space S . This space can either be used for housing, or it remains as undeveloped land. The total housing stock in the city at time t is H_t and the remaining undeveloped land is L_t , so $S = H_t + L_t$ for all t .

A continuum of real estate developers invest in land and construct housing from the land at a cost of K per unit of housing. The aggregate supply of new housing is ΔH_t . Construction is instantaneous, and housing does not depreciate: $H_t = \Delta H_t + H_{t-1}$. Construction is also irreversible: $\Delta H_t \geq 0$. Both housing and land are continuous variables, and one unit of housing requires one unit of land.

The developers rent out land on spot markets at a price of r_t^l . Rental demand for undeveloped land comes from firms, such as farms, that use the city's land as an input. These firms buy their inputs and sell their products on the global market. Therefore, their aggregate demand for land depends only on r_t^l and not on any other local market conditions. This aggregate rental demand curve is $D^l(r_t^l)$, where $D^l(\cdot)$ is decreasing positive function such that $D^l(0) \geq S$.

The profit flow of a developer j at time t is

$$\pi_{j,t} = \underbrace{r_t^l L_{j,t} + p_t^l (L_{j,t-1} - L_{j,t})}_{\text{development profit}} + \underbrace{(p_t^h - p_t^l - K) \Delta H_{j,t}}_{\text{homebuilding profit}}, \quad (2.1)$$

where p_t^h is the price of housing and p_t^l is the price of land. The real estate development industry faces no entry costs, so the industry is perfectly competitive. Because homebuilding is instantaneous and does not depend on prior land investments, profits from this line of

business must be zero due to perfect competition. We denote the aggregate homebuilding profit by $\pi_t^{hb} = (p_t^h - p_t^l - K)\Delta H_t$.

Each developer begins with a land endowment and issues equity to finance its land investments. It maximizes its expected net present value of profits $E_j \sum_{t=0}^{\infty} \beta^t \pi_{j,t}$. The operator E_j reflects firm j 's expectation of future land prices. Firm-specific beliefs represent the beliefs of the firm's CEO, who owns equity, cannot be fired, and decides the firm's land investments. The number of each developer's equity shares equals the amount of land it holds, and each developer pays out its land rents as dividends. The market price of developer equity therefore equals the market price p_t^l of land.

Individual Housing Demand. A population of residents live in the city and hold its housing. These residents receive direct utility from consuming housing. Lower-case h denotes the flow consumption of housing, whereas upper-case H denotes the asset holding. Flow utility from housing depends on whether housing is consumed through owner-occupancy or under a rental contract. Residents also derive utility from non-local consumption c . Each resident i maximizes the expected present value of utility, given by

$$E_i \sum_{t=0}^{\infty} \beta^t u_i(c_t, h_t^{own}, h_t^{rent}),$$

where β is the common discount factor.

Flow utility $u_i(\cdot, \cdot, \cdot)$ has three properties. First, it is separable and linear in non-real estate consumption c . This quasi-linearity eliminates risk aversion and hedging motives. Second, owner-occupied and rented housing are substitutes, and residents vary in which type of contract they prefer and to what degree. Substitutability of owner-occupied and rented housing fully sorts residents between the two types of contracts; no resident consumes both types of housing simultaneously. Finally, residents face diminishing marginal utility of owner-occupied housing. This property leads homeownership to be dispersed among residents in equilibrium.

The utility specification we adopt that features these three properties is

$$u_i(c, h^{own}, h^{rent}) = c + v(a_i h^{own} + h^{rent}) \quad (2.2)$$

where $a_i > 0$ is resident i 's preference for owner-occupancy, and $v(\cdot)$ is an increasing, concave function for which $\lim_{h \rightarrow 0} v'(h) = \infty$. The distribution of the owner-occupancy preference parameter a_i across residents is given by a continuously differentiable cumulative distribution function F_a , which is stable over time. Owner-occupancy utility is unbounded: dF_a has full support on \mathbb{R}^+ . The functional form of the owner-occupancy preference in (2.2) results from a moral hazard problem we describe in the Appendix.

Resident Optimization. Residents hold three assets classes: bonds B , housing H , and developer equity Q . Global capital markets external to the city determine the gross interest rate on bonds, which is $R_t = 1/\beta$, where β is the common discount factor. Residents may borrow or lend at this rate by buying or selling these bonds in unlimited quantities.

In contrast, housing and developer equity are traded within the city, and equilibrium conditions determine their prices p_t^l and p_t^h . Homeowners earn income by renting out the housing they own in excess of what they consume. The spot rental price for housing is r_t^h ; landlord revenue is therefore $r_t^h(H_{i,t} - h_{i,t}^{own} - h_{i,t}^{rent})$. Shorting housing is impossible, but residents can short developer equity. Doing so is costly. Residents incur a convex cost $k_s(Q)$ to short Q units of developer stock, where $k_s(0) = 0$ and $k'_s, k''_s > 0$. These costs reflect fees paid to borrow stock, as well as time spent locating available stock (D'Avolio, 2002).

Short-sale constraints in the housing market result from a lack of asset interchangeability. Although housing is homogeneous in the model, empirical housing markets involve large variation in characteristics across houses. This variation in characteristics makes it essentially impossible to cover a short. Unlike in the housing market, asset interchangeability holds in the equity market, where all of a firm's shares are equivalent.

The Bellman equation representing the resident optimization problem is

$$V(B_{i,t-1}, H_{i,t-1}, Q_{i,t-1}) = \max_{\substack{B_{i,t}, H_{i,t}, Q_{i,t} \\ c_{i,t}, h_{i,t}^{own}, h_{i,t}^{rent}}} \underbrace{c_{i,t} + v(a_i h_{i,t}^{own} + h_{i,t}^{rent})}_{\text{flow utility}} + \underbrace{\beta \mathbf{E}_{i,t} V(B_{i,t}, H_{i,t}, Q_{i,t})}_{\text{continuation value}}, \quad (2.3)$$

where the maximization is subject to the short-sale constraint

$$0 \leq H_{i,t},$$

the ownership constraint

$$h_{i,t}^{own} \leq H_{i,t},$$

and the budget constraint

$$\underbrace{R_t B_{i,t-1} - B_{i,t}}_{\text{borrowing costs}} + \underbrace{c_{i,t}}_{\text{consumption}} \leq \underbrace{p_t^h (H_{i,t-1} - H_{i,t})}_{\text{housing returns}} + \underbrace{r_t^h (H_{i,t} - h_{i,t}^{own} - h_{i,t}^{rent})}_{\text{housing rental income}} \\ + \underbrace{p_t^l (Q_{i,t-1} - Q_{i,t})}_{\text{equity returns}} + \underbrace{r_t^l Q_{i,t}}_{\text{dividends}} - \underbrace{\max(0, k_s(-Q_{i,t}))}_{\text{shorting costs}}.$$

Aggregate Demand and Beliefs. Aggregate demand to live in the city equals the number of residents N_t . This aggregate demand consists of a shock and a trend:

$$\underbrace{\log N_t}_{\text{demand}} = \underbrace{z_t}_{\text{shock}} + \underbrace{\log \bar{N}_t}_{\text{trend}}.$$

The trend component grows at a constant positive rate g : for all $t > 0$,

$$\log \bar{N}_t = g + \log \bar{N}_{t-1}.$$

The shocks z_t have a common factor x . The dependence of the time- t shock on the common factor x is μ_t , so that

$$z_t = \mu_t x.$$

Without loss of generality, $\mu_0 = 1$: the time 0 shock z_0 equals the common factor x . We denote $\boldsymbol{\mu} = \{\mu_t\}_{t \geq 0}$.

At time 0, residents observe the following information: the current and future values of

trend demand \bar{N}_t , the trend growth rate g , the current demand N_0 , the current shock z_0 , and the common factor x of the future shocks. They do *not* observe μ , the data needed to extrapolate the factor x to future shocks. Residents learn the true value of the entire vector μ at time $t = 1$. The resolution of uncertainty at time $t = 1$ is common knowledge at $t = 0$.

Residents agree to disagree about the true value of μ . At time 0, resident i 's subjective prior of μ is given by F_i , an integrable probability measure on the compact space M of all possible values of μ . These priors vary across residents. The resulting subjective expected value of each μ_t is $\mu_{i,t} = \int_M \mu_t dF_i$, and the vector of resident i 's subjective expected values of each μ_t is $\mu_i = \{\mu_{i,t}\}_{t \geq 0}$. The subjective expected value μ_i uniquely determines the prior F_i . The distribution of μ_i itself across residents admits an integrable probability distribution F_μ on M , which is independent from the distribution F_a of owner-occupancy preferences. The CEOs of the development firms are city residents, so their beliefs are drawn from the same distribution F_μ .

Resident disagreement reflects the unprecedented nature of the demand shock z . As argued by Morris (1996), this heterogeneous prior assumption is most appropriate when investors face an unprecedented situation in which they have not yet had a chance to collect information and engage in rational updating. The events surrounding housing booms are precisely these types of situations. Glaeser (2013) meticulously shows that in each of the historical booms he analyzes, reasonable investors could agree to disagree about future real estate prices. In the case of the U.S. housing boom between 2000 and 2006, we follow Mian and Sufi (2009) in thinking of the shock as the arrival of new securitization technologies that expanded credit to low-income borrowers. The initial shock to housing demand is x , and μ represents the degree to which this expansion of credit in 2000-2006 persists after 2006.

Equilibrium. Equilibrium consists of time-series vectors of prices $\mathbf{p}^L(\mu)$, $\mathbf{p}^H(\mu)$, $\mathbf{r}^l(\mu)$, $\mathbf{r}^h(\mu)$ and quantities $\mathbf{L}(\mu)$, $\mathbf{H}(\mu)$ that depend on the realized value of μ . These pricing and quantity functions constitute an equilibrium when housing, land, and equity markets clear while residents and developers maximize utility and profits:

Consider pricing functions $\mathbf{p}^h(\mu)$, $\mathbf{p}^l(\mu)$, $\mathbf{r}^h(\mu)$, $\mathbf{r}^l(\mu)$ and quantity functions $\mathbf{H}(\mu)$, $\mathbf{L}(\mu)$. Let

$H_{i,t}^*$, $Q_{i,t}^*$, $(h_{i,t}^{own})^*$, and $(h_{i,t}^{rent})^*$ be resident i 's solutions to the Bellman equation (2.3) given his owner-occupancy preference a_i , his beliefs μ_i , and these pricing functions. Let $L_{j,t}^*$ be developer j 's land holdings that maximize expected net present value of profits in equation (2.1), given the pricing functions; L_t^* is the sum of these land holdings across developers. The pricing and quantity functions constitute a recursive competitive equilibrium if at each time t :

1. The sum of undeveloped land and housing equals the city's endowment of open space:

$$S = L_t(\mu) + H_t(\mu).$$

2. Flow demand for land equals investment demand from developers, which equals the resident demand for their equity:

$$L_t(\mu) = L_t^* = D^l(r_t^l(\mu)) = \int_0^\infty \int_M Q_{i,t}^* dF_\mu dF_a.$$

3. Resident stock and flow demand for housing clear:

$$H_t(\mu) = N_t(\mu) \int_0^\infty \int_M H_{i,t}^* dF_\mu dF_a = N_t(\mu) \int_0^\infty \int_M ((h_{i,t}^{own})^* + (h_{i,t}^{rent})^*) dF_\mu dF_a.$$

4. Construction maximizes developer profits:

$$H_t(\mu) - H_{t-1}(\mu) \in \arg \max_{\Delta H_t} \pi_t^{hb}.$$

5. Developer profit from homebuilding is zero:

$$\max_{\Delta H_t} \pi_t^{hb} = 0.$$

Elasticity of Housing Supply. The housing supply curve is the city's open space S less the rental demand for land $D^l(r_t^l)$. We denote the elasticity of this supply curve with respect to housing rents r_t^h by ϵ_t^S . The supply elasticity determines the construction response to the shocks $\{z_t\}$. It will also serve as a sufficient statistic for the extent to which land speculation affects house prices. This section describes the supply elasticity ϵ_t^S along the city's trend growth path, which obtains when $x = 0$.

The relationship between land rents r_t^l and house rents r_t^h allows us to calculate this elasticity. Because trend growth $g > 0$, new residents perpetually arrive to the city, and developers build new houses each period. Perpetual construction ties together land and house prices. In particular, as developers must be indifferent between building today or tomorrow, house rents equal land rents plus flow construction costs:

$$r_t^h = r_t^l + (1 - \beta)K.$$

The supply of housing is open space net of flow land demand: $S - D^l(r_t^h - (1 - \beta)K)$. The elasticity of housing supply is thus $\epsilon_t^S \equiv -r_t^h(D^l)' / (S - D^l)$. When the flow land demand D^l features a constant elasticity ϵ^l , the elasticity of housing supply takes on the simple form

$$\epsilon_t^S = \frac{r_t^h}{r_t^h - (1 - \beta)K} \left(\frac{S}{H_t} - 1 \right) \epsilon^l, \quad (2.4)$$

where H_t is the housing stock at time t . The arrival of new residents increases both rents r_t^h and the level of development H_t/S . The supply elasticity given in (2.4) unambiguously falls (see Appendix for proof):

Lemma 3. *Define housing supply to be the residual of the city's open space S minus the flow demand for land: $S - D^l$. The elasticity ϵ_t^S of housing supply with respect to housing rents r_t^h decreases with the level of city development H_t/S along the city's trend growth path.*

2.4 Supply-Side Speculation

At time 0, residents disagree about the future path of housing demand. Speculative trading behavior results from this disagreement. This section describes how owner-occupancy frictions crowd speculators out of owner-occupied housing and into rental housing and land. Demand and supply elasticities determine how prices aggregate the beliefs owner-occupants and of optimistic speculators holding rental housing and land.

2.4.1 Land Speculation and Dispersed Homeownership

We first consider the developer decision to hold land at time 0. Developer j 's first-order condition on its land-holding $L_{j,0}$ is

$$\underbrace{1/\beta}_{\text{risk-free rate}} \geq \underbrace{\mathbf{E}_j p_1^l / (p_0^l - r_0^l)}_{\text{expected land return}},$$

with equality if and only if $L_{j,0} > 0$. A developer invests in land if and only if it expects land to return the risk-free rate. At time 0, developers disagree about this expected return on land because they disagree about the future path of housing demand. The developers that expect the highest returns invest in land, while all other developers sell to these optimistic firms and exit the market. We denote the optimistic belief of the developers who invest in land by $\tilde{\mathbf{E}}p_1^l \equiv \max_{\mu_j} \mathbf{E}(p_1^l \mid \mu_j)$.

Optimistic residents finance developer investments in land through purchasing their equity. Less optimistic residents choose to short-sell developer stock. Developer stock allows residents to hold land indirectly: its price is p_0^l and it pays a dividend of r_0^l . Resident i holds this equity only if he agrees with the land valuation of the optimistic developers, in which case $\mathbf{E}_i p_1^l = \tilde{\mathbf{E}}p_1^l$. Otherwise, he shorts the equity, and his first-order condition is

$$k'_s(-Q_{i,0}^*) = \beta(\tilde{\mathbf{E}}p_1^l - \mathbf{E}_i p_1^l).$$

Disagreement increases the short interest in the equity of the developers holding the land. Without disagreement, $\tilde{\mathbf{E}}p_1^l = \mathbf{E}_i p_1^l$ for all residents, so no one shorts.

Only the most optimistic residents hold housing as landlords. A resident is a *landlord* if he owns more housing than he consumes through owner-occupancy: $H_i > h_i^{own}$. The first-order condition of the Bellman equation (2.3) with respect to $H_{i,0}$ when it is in excess of $h_{i,0}^{own}$ is

$$\underbrace{1/\beta}_{\text{risk-free rate}} \geq \underbrace{\mathbf{E}_i p_1^h / (p_0^h - r_0^h)}_{\text{expected housing return}}, \quad (2.5)$$

with equality if and only if $H_{i,0} > h_{i,0}^{own}$. Only the most optimistic residents invest in rental housing, just as only the most optimistic developers invest in land. Land and rental

housing share this fundamental property. During periods of uncertainty, the most optimistic investors are the sole holders of these asset classes.

Owner-occupancy utility crowds these optimistic investors out of owner-occupied housing, which remains dispersed among residents of all beliefs. The decision to own or rent emerges from the first-order conditions of the Bellman equation (2.3) with respect to $h_{i,0}^{own}$ and $h_{i,0}^{rent}$. We express these equations jointly as

$$\underbrace{v'(a_i(h_{i,0}^{own})^* + (h_{i,0}^{rent})^*)}_{\text{marginal utility of housing}} = \min \left(\underbrace{a_i^{-1}(p_0^h - \beta \mathbf{E}_i p_1^h)}_{\text{owning}}, \underbrace{r_0^h}_{\text{renting}} \right). \quad (2.6)$$

The left term in the parentheses denotes the expected flow price of marginal utility v' from owning a house; the right term denotes the flow price of renting. A resident owns when the owner-occupancy price is less than the rental price. As long as the owner-occupancy preference a_i is large enough, resident i decides to own even if his belief $\mathbf{E}_i p_1^h$ is quite pessimistic. Homeownership remains dispersed among residents of all beliefs.

We gain additional intuition about the own-rent margin by substituting (2.5) into (2.6). We denote the most optimistic belief about future house prices, the one held by landlords investing in rental housing, by $\tilde{\mathbf{E}} p_1^h \equiv \max_{\mu_i} \mathbf{E}(p_1^h | \mu_i)$. The decision to own rather than rent reduces to

$$a_i \geq 1 + \frac{\beta(\tilde{\mathbf{E}} p_1^h - \mathbf{E}_i p_1^h)}{r_0^h}. \quad (2.7)$$

Without disagreement, a resident owns exactly when he intrinsically prefers owning to renting, so that $a_i \geq 1$. Disagreement sets the bar higher. Some pessimists for whom $a_i \geq 1$ choose to rent because they expect capital losses on owning a home. Other pessimists continue to own because their owner-occupancy utility is high enough to offset the fear of capital losses. Proposition 4 summarizes these results.

Proposition 4. *Owner-occupancy utility crowds speculators out of the owner-occupied housing market and into the land and rental markets. The most optimistic residents—those holding the highest value of $\mathbf{E}_i p_1^h$ —buy up all rental housing and finance optimistic developers who purchase all the land. In contrast, owner-occupied housing remains dispersed among residents of all beliefs.*

Proposition 4 yields two corollaries that match stylized facts presented in Section 2.2. The most optimistic developers buy up all the land. Unless they start out owning all the land, these optimistic developers increase their land positions following the demand shock. They hold this land as an investment rather than for immediate construction.

Implication 1. *The developers who hold land at time 0 increase their aggregate land holdings at time 0. They buy land in excess of their immediate construction needs.*

This implication explains the land-buying activities of large public U.S. homebuilders documented in Figure 2.4(a).

The second corollary concerns short-selling. Residents who disagree with the optimistic valuations of developers short their equity.

Implication 2. *Disagreement increases the short interest of developer equity at time 0.*

Figure 2.4(c) documents the rising short interest in the stocks of U.S. public homebuilders who were taking on large land positions during the boom. This short interest provides direct evidence of disagreement during the boom.

2.4.2 Belief Aggregation

Prices aggregate the heterogeneous beliefs of residents and developers holding housing and land. The real estate market consists of three components: land, rental housing, and owner-occupied housing. The most optimistic residents hold the first two, while the third remains dispersed among owner-occupants. House prices reflect a weighted average of the optimistic belief and the average belief of all owner-occupants. The weight on the optimistic belief is the share of the real estate market consisting of land and rental housing; the weight on the average owner-occupant belief is owner-occupied housing's share of the market.

To derive these results, we take a comparative static of the form $\partial p_0^h / \partial x$. The shock $\mathbf{z} = \mu x$ scales with the common factor x . We differentiate with respect to x at $x = 0$ to explore how prices change as the shocks, and hence the ensuing disagreement, increase.

Our partial derivative holds current demand N_0 constant to isolate the aggregation of future beliefs.

We first use (2.5) to write $p_0^h = r_0^h + \beta \tilde{\mathbf{E}} p_1^h$. The shock increases the optimistic belief $\beta \tilde{\mathbf{E}} p_1^h$, directly increasing prices. It also changes the market rent r_0^h . This rent is determined by the intersection of housing supply and housing demand:

$$\underbrace{S - D^l(r_0^h - (1 - \beta)K)}_{\text{housing supply}} = \underbrace{D_0^h(r_0^h)}_{\text{housing demand}}, \quad (2.8)$$

where

$$\begin{aligned} D_0^h(r_0^h) = & \underbrace{N_0 \int_M \int_0^{1+\beta(\tilde{\mathbf{E}} p_1^h - \mathbf{E}_i p_1^h)/r_0^h} (v')^{-1}(r_0^h) dF_a dF_\mu}_{\text{rental housing}} \\ & + \underbrace{N_0 \int_M \int_{1+\beta(\tilde{\mathbf{E}} p_1^h - \mathbf{E}_i p_1^h)/r_0^h}^\infty a_i^{-1}(v')^{-1} \left(a_i^{-1}(r_0^h + \beta(\tilde{\mathbf{E}} p_1^h - \mathbf{E}_i p_1^h)) \right) dF_a dF_\mu}_{\text{owner-occupied housing}}. \end{aligned} \quad (2.9)$$

The housing demand equation follows from (2.6) and (2.7). We determine the shock's effect on rents by totally differentiating (2.8) with respect to x at $x = 0$, keeping current demand N_0 constant. When the elasticity of housing demand ϵ^D is constant, the resulting comparative static $\partial p_0^h / \partial x$ adopts the simple form given in the following proposition, which we prove in the Appendix.

Proposition 5. *Consider the partial effect of the shock in which current demand N_0 stays constant but future house price expectations $\mathbf{E}_i p_1^h$ change. The change in house prices averages the changes in the optimistic resident belief and the average belief:*

$$\frac{\partial p_0^h}{\partial x} = \frac{\epsilon_0^S + (1 - \chi)\epsilon^D}{\epsilon_0^S + \epsilon^D} \frac{\partial \beta \tilde{\mathbf{E}} p_1^h}{\partial x} + \frac{\chi \epsilon^D}{\epsilon_0^S + \epsilon^D} \frac{\partial \beta \bar{\mathbf{E}} p_1^h}{\partial x}, \quad (2.10)$$

where $\tilde{\mathbf{E}} p_1^h = \max_i \mathbf{E}_i p_1^h$ is the most optimistic belief, $\bar{\mathbf{E}} p_1^h = \int_M \mathbf{E}_i p_1^h dF_\mu$ is the average belief, ϵ_0^S is the elasticity of housing supply at time 0, ϵ^D is the elasticity of housing demand, and $\chi = \int_0^\infty (h_{i,0}^{own})^* dF_a / H_0$ is the share of housing that is owner-occupied when $x = 0$.

The weight on the optimistic belief in Proposition 5 represents the share, on the margin,

of the real estate market owned by speculators. The supply elasticity ϵ_0^S represents land, and $(1 - \chi)\epsilon^D$ represents rental housing. The remaining $\chi\epsilon^D$ represents owner-occupied housing and is the weight on the average owner-occupant belief. The average owner-occupant belief coincides with the unconditional average belief because at $x = 0$, beliefs and tenure choice are independent.

Proposition 5 yields four corollaries on the difference in belief aggregation across cities and neighborhoods. Prices look more optimistic when the weight $(\epsilon_0^S + (1 - \chi)\epsilon^D) / (\epsilon_0^S + \epsilon^D)$ is higher. This ratio is greater when the supply elasticity ϵ_0^S is higher:

Implication 3. *Prices look more optimistic when the housing supply elasticity is higher, i.e. in less developed cities.*

Disagreement reverses the common intuition relating housing supply elasticity and movements in house prices. Elastic supply keeps prices low by allowing construction to respond to demand shocks. But land constitutes a larger share of the real estate market when supply is elastic. Speculators are drawn to the land markets, so elastic supply amplifies the role of speculators in determining prices during periods of disagreement. When supply is perfectly elastic, $\epsilon_0^S = \infty$ and prices reflect only the beliefs of these optimistic speculators:

Implication 4. *When housing supply is perfectly elastic, house prices incorporate only the most optimistic beliefs; they reflect the beliefs of developers and not of owner-occupants.*

Recent research has measured owner-occupant beliefs about the future evolution of house prices.⁸ In cities with elastic housing supply, such as the cities motivating this paper, developer rather than owner-occupant beliefs determine prices. Data on the expectations of homebuilders would supplement the research on owner-occupant beliefs to explain prices in these elastic areas.

Prices aggregate beliefs much better when housing supply is perfectly inelastic ($\epsilon_0^S = 0$) and all housing is owner-occupied ($\chi = 1$). In this case, the price change depends only on the average belief $\bar{E}p_1^h$:

⁸See Landvoigt (2011), Case *et al.* (2012), Burnside *et al.* (2013), Soo (2013), Suher (2013), and Cheng *et al.* (2014).

Implication 5. *When the housing stock is fixed and all housing is owner-occupied, prices reflect the average belief about long-run growth.*

In many settings, such as when investor information equals a signal plus mean zero noise, prices reflect all information when they incorporate the average private belief of all investors. Owner-occupied housing markets with a fixed housing stock display this property, even though short-selling is impossible and residents persistently disagree. These frictions fail to bias prices because homeownership remains dispersed among residents of all beliefs, due to the utility flows that residents derive from housing.

The weight $(\epsilon_0^S + (1 - \chi)\epsilon^D)/(\epsilon_0^S + \epsilon^D)$ on optimistic beliefs is also higher when χ is lower:

Implication 6. *Prices look more optimistic when a greater share of housing is rented.*

Speculators own a greater share of the real estate market when the rental share $1 - \chi$ is higher. Prices bias towards optimistic beliefs in market segments where more of the housing stock is rented.

2.5 The Cross-Section of City Experiences During the Boom

This section explains three puzzling aspects of the U.S. housing boom that occurred between 2000 and 2006. First, large house price booms occurred in elastic cities where new construction historically had kept prices low. Second, the price booms in these elastic areas were as large, if not larger, than those happening in inelastic cities at the same time. Finally, house prices remained flat in other elastic cities that were also rapidly building housing.

To explain these cross-sectional facts, we derive a formula for the total effect of the shock \mathbf{z} on house prices. This formula expresses the house price boom as a function of the city's level of development when the shock occurs. Our analysis up to this point has explored the partial effect of how prices aggregate beliefs $E_i p_1^h$, without specifying how these beliefs are formed. To derive the total effect of the shock, we express the changes in these beliefs in terms of city characteristics and the exogenous demand process. Specifically, we calculate

the partial derivative $\partial \log p_0^h / \partial x$ holding all beliefs fixed at $\mu_i = \mu$, and then use Proposition 5 to derive the total effect of the shock x on house prices. As before, we evaluate derivatives at $x = 0$.⁹

At time 0, each resident expects the shock z_t to raise log-demand at time t by $\mu_t x$. The resulting expected change in rents r_t^h depends on the elasticities of supply and demand at time t :

$$\frac{\partial \log E_0 r_t^h}{\partial x} = \frac{\mu_t}{\epsilon_t^S + \epsilon^D}.$$

This equation follows from price theory. When a demand curve shifts up, a good's price increases by the inverse of the total elasticity of supply and demand. The total effects of the shocks $\{z_t\}$ on the current house price p_0^h follows from aggregating the above equation across all time periods, using the relation $p_0 = E_0 \sum_{t=0}^{\infty} \beta^t r_t^h$:

$$\frac{\partial \log p_0^h}{\partial x} = \frac{\mu}{\tilde{\epsilon}^S + \epsilon^D}. \quad (2.11)$$

The mean persistence of the shock is $\mu = \sum_{t=0}^{\infty} \mu_t \beta^t r_t^h (\epsilon_t^S + \epsilon^D)^{-1} / \sum_{t=0}^{\infty} \beta^t r_t^h (\epsilon_t^S + \epsilon^D)^{-1}$, and $\tilde{\epsilon}^S$ is the *long-run supply elasticity* given by the weighted harmonic mean of future supply elasticities in the city:

$$\tilde{\epsilon}^S \equiv -\epsilon^D + \frac{\sum_{t=0}^{\infty} \beta^t r_t^h}{\sum_{t=0}^{\infty} \beta^t r_t^h (\epsilon_t^S + \epsilon^D)^{-1}}.$$

The higher this long-run supply elasticity, the smaller the shock's impact on current house prices, holding μ fixed.

We now put together the two channels through which the shock changes prices. Equation (2.11) expresses the price change that results when μ is known, and (2.10) describes how prices aggregate residents' heterogeneous beliefs about μ . Proposition 6 states the total effect $d \log p_0^h / dx$, which we formally calculate in the Appendix.

⁹Evaluating derivatives at $x = 0$ describes the model when construction occurs in each period. When x is large enough and the shock z might mean-revert, a construction stop at $t = 1$ is possible and anticipated by residents at $t = 0$. This feature of housing cycles, while important, distracts from our focus on housing booms and how they vary across cities.

Proposition 6. *The total effect of the shock x on current house prices is*

$$\frac{d \log p_0^h}{dx} = \underbrace{\left(\frac{\epsilon_0^S + (1 - \chi)\epsilon^D}{\epsilon_0^S + \epsilon^D} \tilde{\mu} + \frac{\chi\epsilon^D}{\epsilon_0^S + \epsilon^D} \bar{\mu} \right)}_{\text{aggregate belief}} \underbrace{\frac{1}{\tilde{\epsilon}^S + \epsilon^D}}_{\text{pass-through}}, \quad (2.12)$$

where ϵ_0^S is the current elasticity of housing supply, $\tilde{\epsilon}^S$ is the long-run supply elasticity, ϵ^D is the elasticity of housing demand, χ is the share of housing that is owner-occupied, $\tilde{\mu}$ is the mean persistence of the most optimistic belief about μ , and $\bar{\mu}$ is the mean persistence of the average belief.

The first puzzle (2.12) explains is how a city with perfectly elastic housing supply can experience a house price boom. Housing supply is perfectly elastic when $\epsilon_0^S = \infty$. In this case, the house price boom is $\tilde{\mu}x / (\tilde{\epsilon}^S + \epsilon^D)$. This price increase is positive as long as the long-run supply elasticity $\tilde{\epsilon}^S$ is not also infinite.

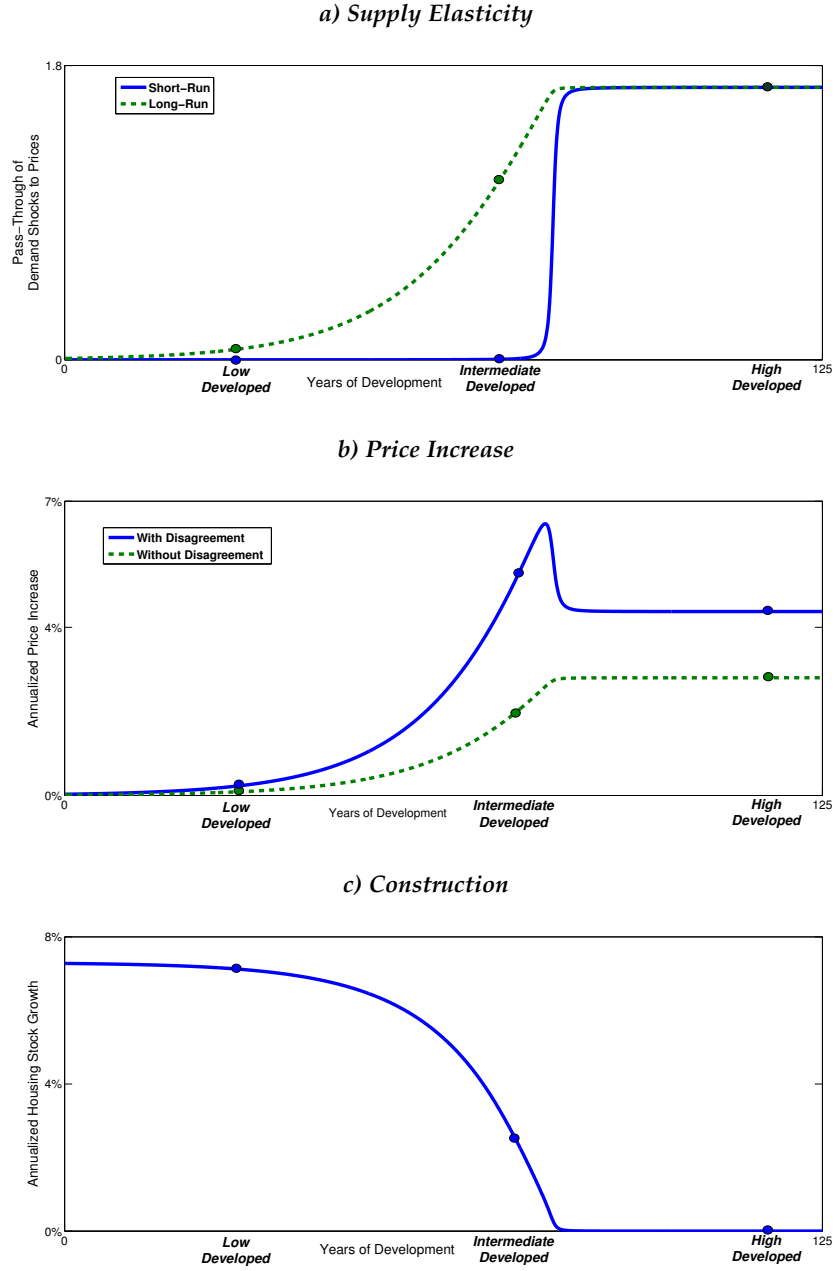
Implication 7. *A house price boom occurs in a city where current housing supply is completely elastic, construction costs are constant, and construction is instantaneous. Supply must be inelastic in the future for such a price boom to occur.*

In the Appendix, we prove that a limiting case exists in which $\epsilon_0^S = \infty$ while $\tilde{\epsilon}^S < \infty$.

A house price boom results from a shock to current demand accompanied by news of future shocks. When supply is inelastic in the long-run, these future shocks raise future rents, and prices rise today to reflect this fact. This price change occurs even if supply is perfectly elastic today, because residents anticipate the near future in which supply will not be able to adjust as easily.

This supply condition—elastic short-run supply, inelastic long-run supply—occurs in cities at an intermediate level of development. Figure 2.5(a) demonstrates the possible combinations of short-run and long-run supply elasticities in a city. We plot the pass-through $1/(\epsilon^S + \epsilon^D)$; a higher pass-through corresponds to a lower elasticity. Lightly developed cities have highly elastic short-run and long-run supply, and heavily developed cities have inelastic short-run and long-run supply. In the intermediate case, current supply is elastic while long-run supply is inelastic, reflecting the near future of constrained supply.

Figure 2.5: Model Simulations For Different Cities



Notes: The parameters we use are $\tilde{\mu} = 1$, $\bar{\mu} = 1$, $x = 0.06$, $g = 0.013$, $\epsilon^D = 1$, $\beta = 0.93^6$, and $\epsilon^I = 1$. We hold the amount of space S fixed and vary the initial trend demand \bar{N}_0 . The x-axis reports annualized trend demand given by $\log \bar{N}_0 / g$. **(a)** Short-run pass-through is $1 / (\epsilon_0^S + \epsilon^D)$; long-run pass-through is $1 / (\bar{\epsilon}^S + \epsilon^D)$. We calculate the rent and housing stock at each level of development using (B.1) in the Appendix, and then calculate the supply elasticities using (2.4). **(b)** Each curve reports the derivative in (2.12) times x , which we calculate using the elasticities shown in panel (a). The “without disagreement” counterfactual uses $\tilde{\mu} = \bar{\mu} = 0.2$ instead of $\tilde{\mu} = 1 > \bar{\mu} = 0.2$. **(c)** We plot the construction equation (B.2) using the elasticities shown in panel (a), as well as rents at each stage of development from (B.1) and prices at each development stage from $p_0 = \sum_{t=0}^{\infty} \beta^t r_t^h$, which we calculate at $x = 0$.

As we discussed in the Introduction, this theoretical supply condition describes the elastic markets that experienced large house price booms between 2000 and 2006. These cities found themselves in a state of arrested development as the boom began in 2000. Although ample land existed for current construction, long-run barriers constrain their future growth.

The second puzzle (2.12) explains why the price booms in these elastic cities were as large as those happening in inelastic cities at the same time. Disagreement amplifies the house price boom the most in exactly these nearly developed elastic cities. The amplification effect of disagreement equals the extent to which optimists bias the price increase given in (2.12). When owner-occupancy frictions are present ($\chi = 1$), the difference between the price boom under disagreement and under the counterfactual in which all residents hold the average belief $\bar{\mu}$ is

$$\frac{\epsilon_0^S}{\epsilon_0^S + \epsilon^D} \frac{\tilde{\mu} - \bar{\mu}}{\tilde{\epsilon}^S + \epsilon^D}.$$

This amplification is largest in nearly developed elastic cities, where ϵ_0^S is large and $\tilde{\epsilon}^S$ is small. Because this amplification increases in ϵ_0^S and decreases in $\tilde{\epsilon}^S$, nearly developed elastic cities provide the *ideal* condition for disagreement to amplify a house price boom. Implication 8, which we prove in the Appendix, states this result formally.

Implication 8. *Disagreement amplifies house price booms most in cities at an intermediate level of development, as long as owner-occupancy frictions are large enough. Define Δ to be difference between the price boom given in (2.12) and the counterfactual in which all residents hold the average belief $\bar{\mu}$. Then there exists $\chi^* < 1$ such that for $\chi^* \leq \chi \leq 1$, Δ is strictly largest at an intermediate level of initial development $\bar{N}_0^* < \infty$.*

Figure 2.5(b) plots the house price boom given by (2.12) across different levels of city development, for both the case of disagreement and the case in which all residents hold the average belief. The amplification effect of disagreement is the difference between the two curves. Optimistic speculators amplify the price boom the most in the intermediate city. Highly elastic short-run supply facilitates speculation in land markets, biasing prices towards their optimistic belief. This bias significantly increases house prices because housing supply

is constrained in the near future. The optimism bias is smaller in the highly developed city. As a result, price increases in intermediate cities are as large as the price boom in the highly developed areas.

In fact, the price boom in some intermediate cities can *exceed* that in the highly developed cities. The parameters we use in Figure 2.5(b) generate an example of this phenomenon. This surprising result reverses the conclusion of standard models of housing cycles, in which the most constrained areas always experience the largest price increases. This reversal occurs as long as owner-occupancy frictions are high and the extent of disagreement is sufficiently large:

Implication 9. *If disagreement and owner-occupancy frictions are large enough, then the largest house price boom occurs in a city at an intermediate level of development. There exists $\chi^* < 1$ and $\delta > 0$ such that for $\chi^* \leq \chi \leq 1$ and $\tilde{\mu} - \bar{\mu} \geq \delta$, the price boom $d \log p_0^h / dx$ is strictly largest at an intermediate level of development $\bar{N}_0^* < \infty$.*

Our model has succeeded in explaining the large house price booms in the elastic cities without arguing that these cities experienced abnormally large housing demand shocks. These markets experienced some of the largest house price booms in the country because of their supply conditions, not in spite of them.

The final puzzle explained by (2.12) is why large house price booms occurred in some elastic cities but not in others. Elastic cities are those for which $\epsilon_0^S \approx \infty$. As shown in Figure 2.5(a), these cities differ in their long-run supply elasticity $\tilde{\epsilon}^S$. When $\tilde{\epsilon}^S = \infty$, the house price boom $d \log p_0^h / dx = 0$. Prices remain flat because construction can freely respond to demand shocks now and for the foreseeable future. House prices increase in elastic cities if and only the development constraint will make construction difficult in the near future.

The elastic American cities which experienced stable house prices between 2000 and 2006 possess characteristics that lead long-run supply to be elastic. These cities, located in Texas and other central American areas, are characterized by flat geography, a lack of future regulation, and homogeneous sprawl (Glaeser and Kahn, 2004; Glaeser and Kohlhase, 2004; Burchfield *et al.*, 2006; Glaeser *et al.*, 2008; Saiz, 2010). These conditions allow the cities to

expand indefinitely, leading S to be infinite or very high. Unlimited land leads the elasticity of supply to remain infinite forever, according to (2.4).

The level of house prices before the shock identifies the difference between the elastic cities that can expand indefinitely and elastic cities that face constraints in the near future. House prices increase with development. Therefore, the elastic cities nearing their development constraints should have higher house prices before the shock than the other elastic cities. The following implication summarizes these results.

Implication 10. *Consider two cities that experience the same demand shock and in which current housing supply is perfectly elastic ($\epsilon_0^S = \infty$). House prices rise more in the city in which the long-run supply elasticity $\tilde{\epsilon}^S$ is lower. Before the shock occurs, a greater share of the land in this city is already developed, and the level of house prices is higher.*

In practice, calculating a metro-area house price level is difficult because characteristics such as construction costs vary widely within and across metro areas, although valiant attempts have been made (Glaeser and Gyourko, 2005; Davis and Heathcote, 2007; Nichols *et al.*, 2010). With the appropriate data, we would be able to distinguish the low-developed from the medium-developed cities.

We have used (2.12) to explain the large house price increases in certain elastic housing markets in the U.S. between 2000 and 2006. An additional salient feature of these booms is that they coincided with rapid construction. As we document in Section 2.2, these cities experienced some of the most intense permitting activity in the nation during this period. Our model captures this phenomenon. Figure 2.5(c) plots the construction response to the shock in different cities. In cities where current housing supply is elastic, new construction accommodates the shock. The elastic cities include both the lightly developed and intermediate developed areas.

2.6 Variation in House Price Booms Within Cities

The model also makes predictions on the variation in house price increases within a given city. Optimistic speculators hold rental housing, just as they hold land. Prices appear more optimistic, and hence house price booms are larger, in market segments where a greater share of housing is rented.

This result emerges from (2.12). Recall that χ is the share of the housing stock that is owner-occupied rather than rented when $x = 0$. It is a sufficient statistic for the distribution F_a of owner-occupancy utility. When χ is larger, the price increase $d \log p_0^h / dx$ is smaller:

$$\frac{\partial}{\partial \chi} \frac{d \log p_0^h}{dx} = -\frac{\epsilon^D}{\epsilon_0^S + \epsilon^D} \frac{\tilde{\mu} - \bar{\mu}}{\tilde{\epsilon}^S + \epsilon^D} < 0.$$

This derivative is negative because the optimistic belief $\tilde{\mu}$ exceeds the average belief $\bar{\mu}$.

A city's housing market consists of a number of market segments, which are subsets of the housing market that attract distinct populations of residents. Because they attract distinct populations, we can analyze them using (2.12), which was formulated at the city-level. All else equal, housing submarkets in which χ is higher experience smaller house price booms:

Implication 11. *Suppose market segments within a city differ only in χ , the relative share of renters versus owner-occupants they attract: the shock x and the short-run and long-run supply elasticities ϵ_0^S and $\tilde{\epsilon}^S$ are constant within a city. Then house price booms are smaller in market segments where χ is larger.*

2.6.1 Location

We first consider variation in χ across neighborhoods. Neighborhoods provide an example of market segments because they differ in the amenities they offer. For instance, some areas offer proximity to restaurants and nightlife; others are characterized by access to good public schools. These amenities appeal differentially to different populations of residents. Variation in amenities hence leads χ to vary across space. Neighborhoods whose amenities appeal relatively more to owner-occupants (high a residents) than to renters (low a residents)

are characterized by a higher value of χ .

Consistent with Implication 11, house prices increased more between 2000 and 2006 in neighborhoods where χ was higher in 2000. We obtain ZIP-level data on χ from the U.S. Census, which reports the share of occupied housing that is owner-occupied, as opposed to rented, in each ZIP code in 2000. The fraction χ varies considerably within cities. Its national mean is 0.71 and standard deviation is 0.17, while the R^2 of regressing χ on city fixed-effects is only 0.12. We calculate the real increase in house prices from 2000 to 2006 using Zillow.com's ZIP-level house price indices. We regress this price increase on χ and city fixed-effects, and find a negative and highly significant coefficient of -0.10 (0.026), where the standard error is clustered at the city level.

However, this negative relationship between χ and price increases may not be causal. Housing demand shocks in this boom were larger in neighborhoods with a lower value of χ . The housing boom resulted from an expansion of credit to low-income households (Mian and Sufi, 2009; Landvoigt *et al.*, 2013), and ZIP-level income strongly covaries with χ .¹⁰

The appeal of χ is that it predicts price increases in any housing boom in which there is disagreement about future fundamentals. In general, χ predicts price increases because it is negatively correlated with speculation, not because it is correlated with demand shocks. Empirical work can test Implication 11 by examining housing booms in which the shocks are independent from χ .

2.6.2 Structure Type

The second approach to measuring χ is to exploit variation across different types of housing structures. According to the U.S. Census, 87% of occupied detached single-family houses in 2000 were owner-occupied rather than rented. In contrast, only 14% of occupied multifamily housing was owner-occupied. According to Implication 11, the enormous difference in χ between these two types of housing causes a larger price boom in multifamily housing, all

¹⁰The IRS reports the median adjusted gross income at the ZIP level. We take out city-level means, and the resulting correlation with χ is 0.40.

else equal.

This result squares with accounts of heightened investment activity in multifamily housing during the boom.¹¹ For instance, a consortium of investors—including the Church of England and California’s pension fund CalPERS—purchased Stuyvesant Town & Peter Cooper Village, Manhattan’s largest apartment complex, for a record price of \$5.4 billion in 2006. Their investment went into foreclosure in 2010 as the price of this complex sharply fell (Segel *et al.*, 2011). Multifamily housing attracts speculators because it is easier to rent out than single-family housing. Optimistic speculators bid up multifamily house prices and cause large price booms in this market segment during periods of uncertainty.

2.7 Conclusion

In this paper, we argue that speculation explains an important part of housing cycles. Speculation amplifies house price booms by biasing prices toward optimistic valuations. We document the central importance of land price increases for explaining the U.S. house price boom between 2000 and 2006. These land price increases resulted from speculation directly in the land market. Consistent with this theory, homebuilders significantly increased their land investments during the boom and then suffered large capital losses during the bust. Many investors disagreed with this optimistic behavior and short-sold homebuilder equity as the homebuilders were purchasing land.

Our emphasis on speculation allows us to explain aspects of the boom that are at odds with existing theories of house prices. Many of the largest price increases occurred in cities that were able to build new houses quickly. This fact poses a problem for theories that stress inelastic housing supply as the source of house price booms. But it sits well with our theory, which instead emphasizes speculation. Undeveloped land facilitates speculation due to rental frictions in the housing market. In our model, large price booms occur in elastic cities facing a development barrier in the near future—cities in arrested development.

¹¹Bayer *et al.* (2013) develop a method to identify speculators in the data. A relevant extension of their work would be to look at the types of housing speculators invest in.

Our approach also makes some new predictions. Price booms are larger in submarkets within a city where a greater share of housing is rented. Although we presented some evidence for this prediction, further empirical work is needed to test it more carefully.

In all, we have presented a different but complementary story of the sources of housing cycles than the literature has offered. Our explanation explains several puzzles and suggests new directions for empirical research.

Chapter 3

Taxation and the Allocation of Talent¹

If we don't have an economy built on bubbles and financial speculation, our best and brightest won't all gravitate towards careers in banking and finance. Because if we want an economy that's built to last, we need more of those young people in science and engineering. This country should not be known for bad debt and phony profits. We should be known for creating and selling products all around the world...

- President Barack Obama, Speech at Osawatomie High School, December 6, 2011

3.1 Introduction

The allocation of talented individuals across professions varies widely across time and space. For example, according to data collected by Goldin *et al.* (2013), more than twice as many male Harvard alumni from the 1969-1972 cohorts pursued careers in both academia and non-financial management as pursued careers in finance. Twenty years later, careers in finance were fifty percent more common than in academia and comparable with those in non-financial management. If private product is anywhere near social product, these talented individuals constitute a large fraction of many societies' human capital: in the United States, for example, just under half of all income is generated by the top 10% of income earners and nearly a fifth is generated by the top 1% (Atkinson *et al.*, 2011).

¹This chapter is co-authored with Benjamin B. Lockwood and E. Glen Weyl.

Furthermore if, as Baumol (1990) and Murphy *et al.* (1991) argue, different professions have different ratios of social to private product (viz. some have negative and others positive externalities) then these differences in talent allocation across societies may have important implications for aggregate welfare. Recent evidence strongly suggests that these externalities not only exist but are large (Murphy and Topel, 2006; Chetty *et al.*, 2013a,b; French, 2008; Piketty *et al.*, 2014). In this paper we argue that non-linear income taxation is a powerful tool affecting the allocation of talent and therefore an important benefit of progressive taxation is increasing aggregate income rather than simply redistributing it.

Our argument is that in selecting an industry, talented individuals face a trade-off between pursuing a “calling” that offers them high non-pecuniary benefits and choosing a career that offers better remuneration. Higher marginal tax rates between the income earned in the lower-paying and higher-paying career make the latter relatively less attractive by narrowing the material sacrifice associated with following passion and prestige. Furthermore, the large shifts of individuals across professions in response to relative wage changes, as well as other evidence from the literature (Lavy and Abramitzky, Forthcoming), suggest (Section 3.3) that the elasticity of such career switches across professions is much larger (1-6) than the intensive elasticity of work within a profession. To the extent that, as suggested by the evidence discussed above, better-paying professions are also more likely to generate negative (less likely to generate positive) externalities, raising marginal tax rates thus has an important pure efficiency benefit.

This framework contrasts with the classical approach pioneered by Vickrey (1945) and Mirrlees (1971), in which labor supply elasticities are the central determinant of optimal tax progressivity. In this respect our approach may be better suited to address contemporary political debates over taxation, which focus more on whether the rich are job creators or “robber barons” than on empirical elasticity estimates. Indeed, we show that externalities calibrated to reflect the views of groups at opposite ends of this debate—the Tea Party and Occupy Wall Street movement—generate optimal income tax schedules that resemble the policies advocated by each group.

In the interest of informing these debates, we turn to the economics literature for estimates of industry-specific externalities. In fact the literature suggests that these externalities are huge and hugely heterogeneous: for each dollar earned privately in research, Murphy and Topel (2006)'s estimates suggest that at least \$5 of positive spill-overs are generated and possibly many more, while French (2008) and Philippon (2013)'s estimates of waste in the financial sector suggest every dollar earned there is accompanied by a 60 cent negative spill-over. As a result our baseline calibration suggests that the allocation of talent has an enormous impact on social welfare and thus on optimal taxation:

- The reallocation of talent between the early 70's and early 90's in the Goldin *et al.* data is sufficient to account for all of the increase in the top 5% income share from its trough mid-century to its peak in 2007 or for almost half of the reduction in growth between the 1948-1973 and the 1982-2007 periods.
- The Reagan (1980 to 1990) tax reforms led to a 3-4 percentage point reallocation of talented individuals from academia, engineering and teaching to finance and management. As a result, these reforms account for a fifth of the change in top incomes in our data and *reduce* total income (inclusive of externalities) by three-quarters of a percent. The reforms reduced social welfare by one to two percent.
- In a simple model with no intensive elasticity or redistributive motive, efficiency alone justifies nearly confiscatory (80-90%) tax rates and subsidies of more than 100% for joining the middle class, while in our richest calibration, which includes realistic intensive margin elasticities and a redistributive motive, optimal top marginal rates are 70-80%.

While some of these calibrated magnitudes strike us as implausibly large, they follow from a conservative, if inevitably partial, reading of the literature and a simple but natural model. Even if, as we suspect, future research shows that the existing literature overestimates externalities and thus our magnitudes, our analysis emphasizes the importance of refining such estimates and incorporating them into the analysis of optimal income taxation, to

which limited attention has been devoted in existing research. Because the externalities that are our central parameters are so uncertain, we have developed an applet, written jointly with Joshua Bosshardt (<http://taxapplet.appspot.com>), that allows the reader to input her desired parameter values and reads off optimal tax rates for any given views she has about externalities.

Income taxation is admittedly a blunt tool in addressing differences in externalities across professions, as our results show in Subsection 3.5.4. In fact only 5-20% of the first-best welfare gains possible under profession-specific taxation are possible under non-discriminatory income taxation. Nevertheless we believe that analyzing the effects of non-discriminatory income taxation is useful for policy for several reasons. First, occupations are easy to conceal or misrepresent, so a tax code that directly rewards or penalizes specific professions would be difficult to enforce, and externalities measurements would be difficult to update. Under an income tax, however, agents have no incentive to misrepresent their profession, and thus the continued measurement of externalities is incentive compatible. Second, the political economy consequences of allowing profession-specific taxation could be dangerous, unleashing a range of special interest lobbying and propaganda that is unlikely to lead to an efficiency-enhancing equilibrium. Third, considerations of horizontal equity may make differential taxation of different occupations ethically or politically unpalatable. Fourth, profession-specific taxation is simply not on the public agenda, while income tax reform is; thus, given the current second-best situation, we believe economists' views of optimal taxation should be influenced by their goals in allocating talent, not just their views about redistribution. Many of these concerns are analogous to those that led Mirrlees (1971) to focus attention on non-linear income taxation and assume that wages as such were non-contractible even though they seem to be at least partially observable in practice.

Rothschild and Scheuer (2014a,b), and to a lesser extent Philippon (2010), employ theoretical models to explore the effects of the economy on optimal income taxation; we discuss the relationship between these papers more extensively in Subsection 3.4.4. Partly as a result, our goal is not only to highlight the use of income taxation to address externalities

created by professions, but additionally to provide a simple and rigorous framework that addresses the quantitative relevance of such factors and contrast them with the classical insurance-incentives trade-off on which most of the optimal tax literature is based.

Following this introduction, the paper is divided into five sections. Section 3.2 develops our ideas in a simple theoretical model that makes a strong (and likely unrealistic) assumption about substitution patterns across professions which renders the formula for optimal taxation particularly transparent. Section 3.3 uses data on income distributions and findings from the economics literature about professions' externalities to calibrate this model. Section 3.4 discusses some of the implications of these results and their relationship to the literature. Section 3.5 presents a structural model based on a more realistic assumption about substitution patterns at the cost of imposing other assumptions that make the results either more special or less transparent. We show in several ways that these more realistic substitution patterns further strengthen our results from the simple case in Section 3.3 and use the model to quantitatively evaluate the importance of our mechanism. Section 3.6 concludes. Details of our empirical procedures and less instructive proofs are in appendices following the main text.

3.2 Theory

In this section we construct the simplest and most intuitive version of our theory by ruling out any redistributive motive and considering the optimization of a non-linear income tax to sort talent across professions to maximize aggregate income-equivalents in the spirit of Kaldor (1939) and Hicks (1939). We consider a static model intended to capture the long-term, steady-state choices of a population. Because career choices are very sticky, they should be interpreted as being made once and for all, likely at an early stage of life. Thus the "career-switching" we refer to is primarily an elasticity of changes in choices likely made by college students, not changes mid-career for adults.

3.2.1 General model

There are N professions, $p = 1, \dots, N$ and a mass 1 of talented individuals. Each individual i is characterized by a $2N$ -dimensional vector $\theta_i = (\mathbf{w}_i, \boldsymbol{\psi}_i)$, where w_{ip} represents the money wage individual i would earn if she chose profession p and ψ_{ip} represents a parameter characterizing the non-pecuniary or *psychic* income she would receive from working in this profession. These characteristics of individuals are distributed in the population according to a non-atomic and differentiable distribution function f with full support on a convex and open $\Theta \subseteq \mathbb{R}_{++}^{2N}$. Each individual can choose how many hours to work, h , at a utility cost $\phi(h; \psi_{ip})$ where $\frac{\partial \phi}{\partial h} > 0$ (work is costly), $\frac{\partial^2 \phi}{\partial h^2} > 0$ (the costs of work are convex) and $\frac{\partial \phi}{\partial \psi} \leq 0$ (ψ shifts down the costs of working/up the non-pecuniary benefits of work). Wages are assumed linear, so the individual i 's income in profession p given that she works h hours is $w_{ip}h$. We do not constrain hours to lie in a finite range as we interpret hours more broadly as effort and we do not assume ϕ need be positive (work may be enjoyable, on net, but the marginal cost of hours is also positive).

The government must finance a net expenditure of I (which we typically assume to be 0) through the use of a non-linear income tax under which an individual earning total income $y = wh$, regardless of the profession in which she earns this wage, pays a total tax $T(y)$. Thus, just as in Mirrlees (1971), we effectively assume that neither wage nor profession can be verified or that, as suggested by Diamond and Saez (2011), some horizontal equity concern prevents greater discrimination.² Each individual's utility is just the sum of her money and psychic incomes.

Each individual is assumed to have quasi-linear utility in money income and thus to earn net utility

$$w_p h - T(w_p h) - \phi(h; \psi_{ip}) + E$$

when she works h hours in profession p , where E is the net externalities she receives from other individuals, as discussed below, which is uniform across the population and thus

²See our discussion in the introduction for a more detailed justification.

independent of her actions. Given the assumed convexity of ϕ , so long as T is also convex (or not “too concave”—see Mirrlees and Saez (2001) for details) each individual has a unique optimal $h_p^*(w_p, \psi_p)$ to work conditional on being in profession p and will only move h_p^* locally in response to small marginal tax changes.³ We let $u_p^*(w_p, \psi_p)$ be the value of utility at this optimum. Sometimes we write $h_p^*(\theta)$, and similarly for u_p^* , which is interpreted as extracting the relevant components. Individual i chooses to work in the profession where her after-tax income plus her psychic income is highest, at her optimal profession-conditional hours level:

$$p^*(\theta) \equiv \operatorname{argmax}_{p \in 1, \dots, N} w_p h_p^*(w_p, \psi_p) - T(w_p h_p^*(w_p, \psi_p)) - \phi(h_p^*(w_p, \psi_p); \psi_p).^4 \quad (3.1)$$

Each profession has an externality share, e_p with the interpretation that an individual i working in profession p and earning wage w_{ip} generates a net externality on the rest of society, evenly distributed across individuals, of $e_p w_{ip}$ for each hour she works. We thus assume a linear technology here and assume away all general equilibrium effects on wages. While this is consistent with the standard Mirrlees approach, it may cause interpretative problems in industries with non-linear technologies as the ratio of marginal and average products can be mistaken for externalities. While a model with non-linear technology would be needed to disentangle such issues rigorously, we believe that our externality shares should be identified with the share of marginal output in the industry that is external and conjecture that this could be shown, at least in first-order conditions, in a broader model.⁵

³We do not constrain T to be convex and some of our optimal schedules have non-convex regions of T . In these cases, when we allow the hours margin to be flexible, we make explicit our assumptions about hours elasticities, which always rule out non-local movements. Given our focus on non-local income movements driven by career switches we view this as a reasonable simplification. Such simplifications are standard in the literature (Saez, 2001).

⁴When the best profession for an individual is not unique, the tie may be broken in any manner as our assumptions on the distribution of wages and psychic incomes assure that the set of individuals facing such indifferences is of measure 0.

⁵Under this view, it is important to distinguish between the marginal and average externalities of a profession. For example, it may be that inframarginal members of a profession generate large positive externalities while marginal individuals may have zero or negative externalities. Similarly it may be that a profession is governed

The value E received by all individuals is thus the average value of $e_{p^*(\theta)} w_{p^*(\theta)} h_{p^*(\theta)}^*(\theta)$. The planner seeks to maximize total income (money net of externalities and psychic income) in her choice of the tax. In particular, she solves

$$\max_{T(\cdot)} \int_{\Theta} \left[\left(1 + e_{p^*(\theta)} \right) w_{p^*(\theta)} h_{p^*(\theta)}^*(\theta) - \phi \left(h_{p^*(\theta)}^*(\theta); \psi_{p^*(\theta)} \right) \right] f(\theta),$$

subject to the definition of h_p^* and p^* above. To derive optimal taxes, we follow the intuitive perturbation approach to the calculus of variations problem pioneered in economics by Wilson (1993) and in optimal income taxation by Saez (2001). Suppose the planner raises slightly the marginal tax rate at income y , returning the raised revenue uniformly to the population so as to continue to satisfy her budget constraint and otherwise leaving fixed all other parts of the tax system. The redistribution thus induced has no net social value, as the planner seeks only to maximize total social wealth. Thus we can focus entirely on the behavioral responses to the tax rise.

One component of this is the local, *intensive* elasticity through the number of hours individuals choose to work. Given convexity of ϕ , so long as T is not too concave the optimal choice of h will always move locally in response to small changes in the optimal tax rate. In particular, if marginal taxes at income y , $T'(y)$, rises we can trace the impact on the optimal hours for an individual of type θ who is earning income y (because $w_{p^*(\theta)} h_{p^*(\theta)}^*(\theta) = y$), assuming she stays in the same profession, by the implicit function theorem. The first-order condition for $h_{p^*(\theta)}^*$ is

$$w_{p^*(\theta)} \left[1 - T' \left(w_{p^*(\theta)} h_{p^*(\theta)}^*(\theta) \right) \right] = \phi' \left(h_{p^*(\theta)}^*(\theta); \psi_{p^*(\theta)} \right). \quad (3.2)$$

We can now determine the effect of increasing $T'(y)$ by a small amount using the implicit function theorem, letting $\epsilon_p^h(\theta)$ be the (negative) *intensive labor supply elasticity* of h with

by a restrictive cartel (e.g. medicine and law) so that marginal entry into the profession mostly crowds others out of the profession, even though overall the profession generates positive or negative externalities.

respect to the post-tax wage $w(1 - T')$:

$$-w_p \left(1 - \frac{\epsilon_p^h(\theta) h_p^*(\theta) T''(y)}{1 - T'(y)} \right) = - \frac{\epsilon_p^h(\theta) h_p^*(\theta) \phi''(h_p^*(\theta); \psi_p)}{w_p [1 - T'(y)]} \implies$$

$$\epsilon_p^h(\theta) = \frac{w_p^2 [1 - T'(y)]}{h_p^*(\theta) [\phi''(h_p^*(\theta); \psi_p) + w_p^2 T''(y)]}. \quad (3.3)$$

On the other hand, the marginal social value created by an individual working an additional hour is

$$(1 + e_{p^*(\theta)}) w_{p^*(\theta)} - \phi'(h_{p^*(\theta)}^*(\theta); \psi_{p^*(\theta)}) = [e_{p^*(\theta)} + T'(y)] w_{p^*(\theta)},$$

where the equality follows by substituting in the first-order condition for hours, equation (3.2). Intuitively, by the envelope theorem, the net social value created by an additional hour of work is proportional to the private product (wage) multiplied by total externality associated with wages earned, both through the tax externality and the direct externality e .

If marginal tax rates rise at y , an individual of type θ currently earning income y will change her hours by $\frac{\epsilon_{p^*(\theta)}^h(\theta) h_{p^*(\theta)}^*(\theta)}{1 - T'(y)}$ and thus will change social welfare by

$$\frac{[e_{p^*(\theta)} + T'(y)] y \epsilon_{p^*(\theta)}^h(\theta)}{1 - T'(y)}.$$

The set of all individuals earning income y is

$$\Theta(y; T) \equiv \left\{ \theta \in \Theta : w_{p^*(\theta)} h_{p^*(\theta)}^*(\theta) = y \right\}.$$

Where it does not create ambiguity below, we drop the dependence on T . Applying the multidimensional Leibniz rule of Veiga and Weyl (2013), the density (normalizing for the representation in the type-space) of such consumers is

$$f(y) \equiv \int_{S \in \Theta(y)} \frac{f(\theta(S))}{-\frac{\partial \phi}{\partial \psi_{p^*(\theta(S))}}(\theta(S))} dS,$$

where S is a $2T - 1$ -dimensional parameterization of the set $\Theta(y)$. As short hand, we can abbreviate this notation as $\int_{\Theta(y)} f(\theta) d\theta$. Under this notation, the total impact on welfare

through these local changes is then

$$\int_{\Theta(y)} \frac{[e_{p^*(\theta)} + T'(y)] \epsilon_{p^*(\theta)}^h(\theta)}{1 - T'(y)} f(\theta) d\theta =$$

$$\frac{f(y) \left(\left[E[e_{p^*(\theta)} | \Theta(y)] + T'(y) \right] E[\epsilon_{p^*(\theta)}^h(\theta) | \Theta(y)] + \text{Cov}(e_{p^*(\theta)}, \epsilon_{p^*(\theta)}^h(\theta) | \Theta(y)) \right)}{1 - T'(y)},$$

where the expectation and covariance operators are defined as usual, conditional on the relevant sets in the short-hand notation; see Veiga and Weyl (2013) for greater details.

The second component of the behavioral response follows a similar normative logic, but is driven by changes in professions. In particular, let

$$\partial\Theta(y; T) \equiv \left\{ \theta \in \Theta : \exists p, q \in 1, \dots, N : \left(w_p h_p^*(\theta) < y < w_q h_q^*(\theta) \right) \wedge \left(u_p^*(\theta) = u_q^*(\theta) \right) \right\}$$

be the set of *y-career switching* individuals who, under tax system T , are just indifferent between two professions, one of which has an optimal (for that individual) income level above and the other of which has an optimal income level below y . Let $f_S(y) \equiv \int_{\partial\Theta(y)} f(\theta) d\theta$ be the density of such individuals. Raising $T'(y)$ causes all of these individuals to switch from profession q to profession p .⁶ How does this change social wealth created?

In profession q , the individual i generates social wealth $(1 + e_{iq}) w_{iq} h_q^*(\theta_i) - \phi(h_q^*(\theta_i), \psi_{iq})$ while in profession p she generates social wealth $(1 + e_{ip}) w_{ip} h_p^*(\theta_i) - \phi(h_p^*(\theta_i), \psi_{ip})$. From the fact that she is indifferent between the two professions, we know that

$$w_{ip} h_p^*(\theta_i) - T(w_{ip} h_p^*(\theta_i)) - \phi(h_p^*(\theta_i), \psi_{ip}) =$$

$$w_{iq} h_q^*(\theta_i) - T(w_{iq} h_q^*(\theta_i)) - \phi(h_q^*(\theta_i), \psi_{iq}).$$

Thus the change in social wealth created by her switching professions is

$$(1 + e_{ip}) w_{ip} h_p^*(\theta_i) - \phi(h_p^*(\theta_i), \psi_{ip}) - \left[(1 + e_{iq}) w_{iq} h_q^*(\theta_i) - \phi(h_q^*(\theta_i), \psi_{iq}) \right] =$$

⁶Note we can ignore individuals who are triply indifferent between professions as they are of measure zero even within the career switchers.

$$e_{ip}w_{ip}h_p^*(\theta_i) + T(w_{ip}h_p^*(\theta_i)) - e_{iq}w_{iq}h_q^*(\theta_i) - T(w_{iq}h_q^*(\theta_i)) \equiv \Delta T(\theta_i) + \Delta E(\theta_i),$$

that is the change in the sum of her tax payments and externalities. This is exactly the discrete, career-switching analog of the intensive margin change in hours. The total change in social wealth from an increase in the marginal tax rate at y is $E[\Delta T(\theta) + \Delta E(\theta) | \partial\Theta(y)] f_S(y)$.

Socially optimal taxation calls for equating the sum of the intensive and career-switching effects to 0:

Proposition 7. *Assuming all movements in hours are local, optimal taxation requires that for all $y : f_S(w)$ or $f(y)$, $E[\epsilon_{p^*}^h | \Theta(y)] > 0$,*

$$\frac{\left([E[e_{p^*} | \Theta(y)] + T'(y)] E[\epsilon_{p^*}^h | \Theta(y)] + \text{Cov}(e_{p^*}, \epsilon_{p^*}^h | \Theta(y)) \right) f(y)}{1 - T'(y)} + E[\Delta T + \Delta E | \partial\Theta(y)] f_S(y) = 0.$$

While this intuitive result applies quite generally, it provides relatively little guidance on how to map information about distribution of incomes and externality shares to optimal tax rates in the absence of detailed information on which individuals are likeliest to switch professions. To clarify the analysis further, we now consider additional assumptions that may be added to yield an especially simple formula.

3.2.2 A simple case

The first such assumption is that there is no correlation between propensity to switch careers, conditional on an income level, and the externalities generated by an individual earning that income nor between the propensity to switch into a career and the externalities, conditional on the income level switched into.

Assumption 1. *For an individual in a switching set, $\partial\Theta(y)$ for some y , let $p(\theta)$ be the lower-paying of the two professions she is indifferent between and $q(\theta)$ be the higher-paying of the two professions she is indifferent between. For all y, y'*

$$E[e_{p^*} | \Theta(y)] = E[e_p | \partial\Theta(y'), w_p h_p^* = y] = E[e_q | \partial\Theta(y'), w_q h_q^* = y].$$

That is, the average individual considering (in response to a marginal change in the tax rate at y') switching either down to a lower-paying profession or up to a higher-paying profession but currently earning income y on average generates the same externality as the average individual earning that income.

The second assumption is analogous, but for the intensive margin: elasticities are uncorrelated with externalities, conditional on income.

Assumption 2. *For every y ,*

$$\text{Cov} \left(e_{p^*(\theta)}, \epsilon_{p^*(\theta)}^h | \Theta(y) \right) = 0.$$

A simple example in which these assumptions would be satisfied is the case where those exiting and entering any income level are chosen randomly and representatively from the skilled professions at that income level. Alternatively, to illustrate that Assumptions 1 and 2 can be consistent with micro-foundations, consider the following more stylized example:

Example 1. *Suppose that the equilibrium incomes of different professions have disjoint supports.⁷ Then at any income Assumptions 1 and 2 are automatically satisfied as the only individuals earning any given income all have the same externality share. Conditions on primitives generating this would be that all individuals find hours below a certain range to have no cost and above a certain range to be infinitely expensive; call this range $[\underline{h}, \bar{h}]$. The range of incomes generated by individuals in profession p are then $[\underline{w}_p \underline{h}, \bar{h} \bar{w}_p]$ where \underline{w}_p is the lowest value w_p may take on and \bar{w}_p is the highest value it may take on. If $[\underline{w}_p \underline{h}, \bar{h} \bar{w}_p]$ are disjoint for different p then so will income be.*

Although the conditions in this example are not satisfied in the data presented in Section 3.3.1, which has overlap at many levels of the distribution, the example does show that the assumptions can be satisfied and that they may be approximately satisfied if professions are highly segregated by income. In Section 3.5 we provide stylized but natural alternative models in which the results are even more extreme, in terms of the responsiveness of optimal taxation to externality shares, than the results we derive in this section.

⁷We are grateful to Florian Scheuer for suggesting this example.

Under Assumptions 1 and 2, optimal taxation takes a very simple form in both the case of pure career-switching ($\epsilon_{p^*}^h(\theta) \equiv 0 \forall \theta \in \Theta$) and in the case of pure intensive reactions and no career switching ($f_S(y) = 0 \forall y$):

Proposition 8. *Suppose that, for every tax policy T , Assumption 1 holds and that hours are rigid (every individual must work an exogenous number of hours in each profession). Then*

$$T^*(y) = -E[e_{p^*} | \Theta(y; T)] y + T_0,$$

where T_0 is a constant across income levels. That is, up to a lump sum transfer, average tax rates are set at each income to offset the average externality created by individuals earning that income level. This is true if the average externalities are defined at the current equilibrium or at the optimal policy, as the average externalities at each income level remain the same in this case.

Alternatively suppose that, for every tax policy, Assumption 2 holds, careers are rigid and all movement in hours is local. Then

$$T^{*'}(y) = -E[e_{p^*} | \Theta(y; T^*)].$$

That is, marginal tax rates at each income level are set to offset the average externality created by individuals earning that income level given the optimal tax policy.

Intuitively, when there is pure career switching the average externality of each profession must offset the average cost of each profession because when switching across professions it is the average income rather than the marginal earnings that are relevant. On the other hand, when making marginal decisions about work, it is marginal tax rates that are relevant and thus marginal tax rates should be equated to the average externality at a given income level. This intuition is formalized in the proof of this result in Appendix C.8. As we discuss in the next section, the distinction between policies generated by these two regimes is fairly small in many contexts.

Proof. See Appendix C.8. □

We refer to the optimal policy under pure career-switching as “average tax externality

matching” (ATEM) and the optimal policy under pure intensive hours choice “marginal tax externality matching” (MTEM). To see the distinction between these policies, note that under ATEM

$$T^{\star'} = -E[e_{p^*}|\Theta(y; T^{\star})] - y \frac{\partial E[e_{p^*}|\Theta(y; T)]}{\partial y}. \quad (3.4)$$

Suppose that externalities are becoming larger in absolute magnitude as income rises. Then marginal tax rates are more extreme under ATEM than MTEM: if average income-conditional externalities are negative and increasing in size with income and thus marginal tax rates are positive under MTEM they will be even larger under ATEM. If average income-conditional externalities are positive and increasing with income and thus marginal tax rates are negative under MTEM they will be even more negative under ATEM.

Proposition 2 dramatically simplifies the data requirements for determining optimal tax policy. However, it still presents two challenges. First, in the case when there are both career-switching and intensive elasticities, optimal policy is a mix of these two extremes, the mix depending on details of how large the two elasticities are over in the distribution of income. Second, with no career switching, it is the average income-conditional externality *at the optimal policy* rather than at the current equilibrium that is relevant to determine marginal tax rates. This is harder to observe from available data for obvious reasons, though this same challenge appears, and is treated as we do below, in much standard optimal tax work (such as Saez (2001)).

3.3 Calibration

In what follows, we focus primarily on the case of ATEM because we believe that career-switching elasticities are much larger than hours elasticities, at least for the most talented individuals. For example while Saez *et al.* (2012) argue that elasticities of intensive margin labor supply are very low (0 to 0.1 for high incomes), Goldin *et al.*’s evidence suggests career-switching elasticities are much higher. The share of male Harvard alumni who pursued a career in finance, for example, more than tripled from 5% in the 1969-72 cohort

to 15.7% in 1989-92 cohort. Throughout *all* of the analysis that follows we restrict attention within the Harvard data to males, though the IRS data does not distinguish by gender.

Suppose that there have been no changes in inherent preferences for different careers (no labor supply shift) and thus that all of this change arises from shifts in relative wages. Philippon and Reshef (2012) suggests that post-tax wages in the financial sector have increased, from the 1980s to the late 1990s, by somewhere between 50% and 200% depending on how one adjusts for education and expectations; see Appendix C.4 for details. Thus elasticities are somewhere between 1 and 6. Given that these increased wages were likely not fully anticipated and mostly accrued in the 1990s and 2000s and that finance's share of graduates appears, anecdotally, to have greatly increased further since the 89-92 cohort, these elasticities may be underestimates. Even Keane and Rogerson (2012), who favor higher hours elasticities, argue that longer-term elasticities (along dimensions like career choice), are likely to be much larger than intensive hour elasticities. We also show that our results are qualitatively robust to calibrating MTEM.

Computing ATEM and MTEM requires data on two things: income distributions for different industries, and estimates of the externalities (positive or negative) created by each industry. The first is relatively straightforward; in Section 3.3.1 we describe how we use data from IRS tax returns and the salary distributions of Harvard graduates for this estimation.

To calibrate externalities, we adopt two approaches. First, we present estimates of profession-specific externalities from the economics literature. Second, we include two sets of externality assumptions intended to reflect the poles of contemporary debates over income tax progressivity in the United States: the Tea Party and the Occupy Wall Street movement, both to show the sensitivity of the optimal tax schedule to a range of externality assumptions, and to demonstrate the usefulness of this model in formalizing observed political debates.

3.3.1 Income distributions

Due to the thick upper tail of the US income distribution, welfare and optimal taxes depend critically on the allocation of talent among high-earning professions. We therefore restrict our attention to 10 professions which account for nearly all of the top incomes and occupational choices of talented individuals in the United States: law, finance, management, medicine, academia/science, computers/engineering, sales, consulting, teaching (primary and secondary), and arts/entertainment. As our data will show, 88% of the top 1% incomes reported to the IRS come from these professions, and 93% of male students who graduated Harvard College in 1990 work in these occupations in 2005.

We jointly estimate the income distribution within each profession and the share of skilled workers in profession p (denoted $F_p(\cdot)$ and s_p , respectively) using data from two sources. First, we use results reported in Bakija *et al.* (2012), based on IRS tax returns, where we observe the share of the top 1% and 0.1% of the population income distribution (i.e., with earnings greater than \$295,000 and \$1,246,000, respectively) employed in each of our skilled professions (except Teaching, see below) in 2005. We denote these shares b_{1p} and b_{2p} respectively, where p represents profession. Second, we use data from Goldin *et al.* (2013) to construct the empirical CDF of earnings for the Harvard class of 1990 in each all 10 skilled professions. Specifically, we observe a list of incomes $y_0^h < y_1^h < \dots < y_m^h$ where $y_0^h = 0$ and $y_m^h = \infty$ and an $m \times 10$ matrix $(a)_{ip}$ such that a_{ip} is the observed share of the Harvard cohort in occupation p that is earning between y_{i-1}^h and y_i^h .

Our parametric assumption is that incomes within each profession follow a Pareto-lognormal distribution. This distribution resembles a lognormal distribution for low values and a Pareto distribution for high values, and was introduced by Colombi (1990) to model income distributions. Unlike other distributions, it does a good job matching both the top tail and the central mass of the income distribution. The Pareto-lognormal distribution is characterized by three parameters, which we estimate within each profession using maximum likelihood as follows.

Because the IRS data from Bakija *et al.* represent the entire population (whereas the

Harvard data in Goldin et al. represent a subset and we focus attention on males even within this subset), we regard the former as more reliable. Yet because it contains only two data points for each profession, that data is insufficient to calibrate the 3 parameters of Pareto-lognormal distribution, so we select parameters to maximize the likelihood of observing the Harvard data, conditional on matching the IRS data points exactly. Formally, we use $F(\cdot; \mathbf{z})$ to denote the Pareto-lognormal distribution with 3-dimensional parameter \mathbf{z} , and \mathbf{z}_p to denote the true parameter vector characterizing the Pareto-lognormal in profession p . We let s_p denote the share of the skilled population employed in profession p . Thus we select our parameter estimates $(\hat{\mathbf{z}}_p, \hat{s}_p)$ to solve

$$(\hat{\mathbf{z}}_p, \hat{s}_p) = \arg \max_{\mathbf{z}_p, s_p} \sum_{i=1}^m a_{ip} \log \left(F(y_i^h; \mathbf{z}_p) - F(y_{i-1}^h; \mathbf{z}_p) \right)$$

such that $s_p (1 - F(295000; \mathbf{z}_p)) = 0.01b_{1p}$ and $s_p (1 - F(1246000; \mathbf{z}_p)) = 0.001b_{2p}$.

The IRS does not provide data on Teaching (instead lumping teaching in with government). For teaching, we therefore estimate $\mathbf{z}_{\text{teaching}}$ directly from the Harvard data:

$$\hat{\mathbf{z}}_{\text{teaching}} = \arg \max_{\mathbf{z}_{\text{teaching}}} \sum_{i=1}^m a_{i,\text{teaching}} \log \left(F(y_i^h; \mathbf{z}_{\text{teaching}}) - F(y_{i-1}^h; \mathbf{z}_{\text{teaching}}) \right).$$

We estimate the fraction s_{teaching} of talented individuals entering teaching by equating this fraction to the fraction of Harvard graduates in teaching, which is stable across cohorts (1970, 1980, and 1990 graduates) and is approximately 3%.

Figure 3.1 shows the IRS and Harvard data for each profession except finance, as well as the fitted Pareto log-normal distribution. (Note that the IRS data points depend on the shares \hat{s}_p , which are estimated jointly with the Pareto-lognormal parameters.) The fit appears quite quite good.

In Finance, the IRS data and the Harvard data are in direct conflict; the Harvard data shows a far richer upper tail of the income distribution than does the IRS data. This raises some concern that our calibration may underestimate the representation of Finance at very high incomes, a hypothesis consistent with the strong representation of finance that Kaplan and Rauh (2010) find at high incomes. To account for this possibility, we also compute

Figure 3.1: *Income distributions fitted to IRS and Harvard data in 9 industries*

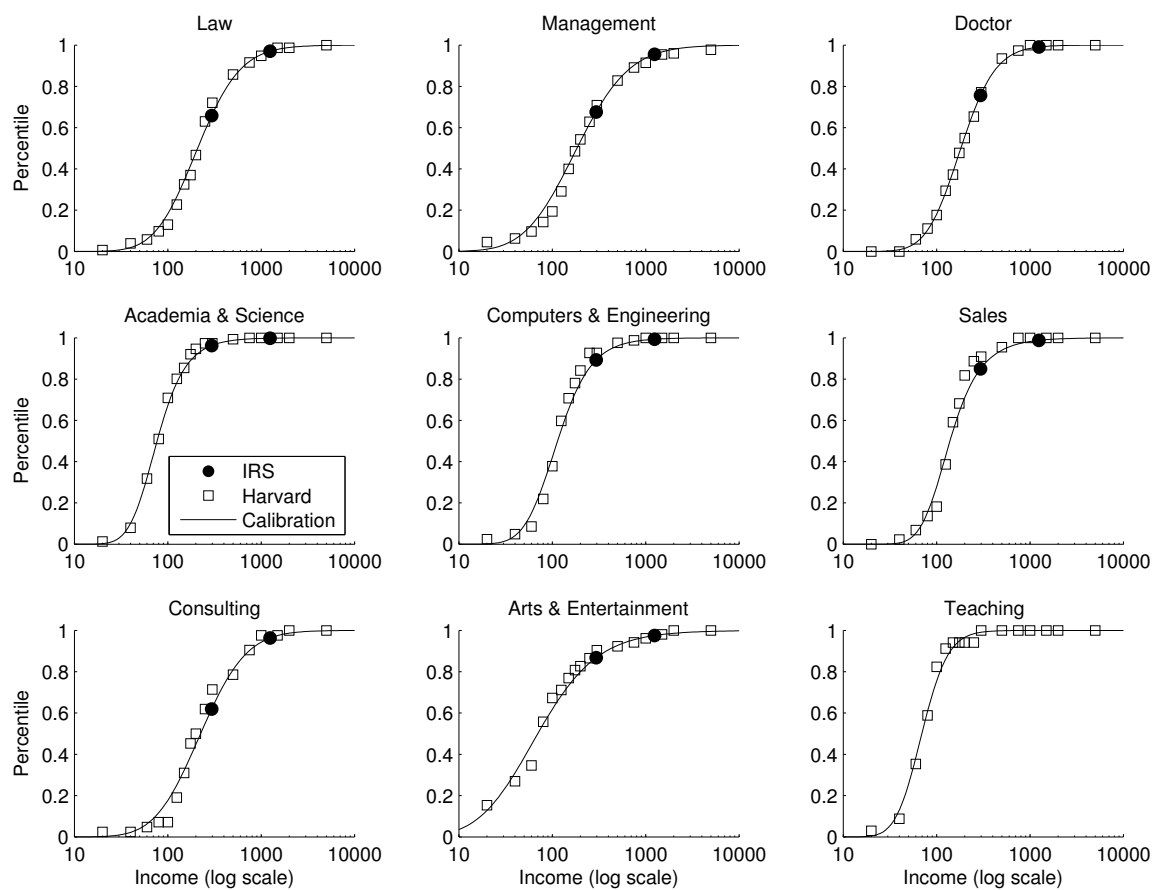
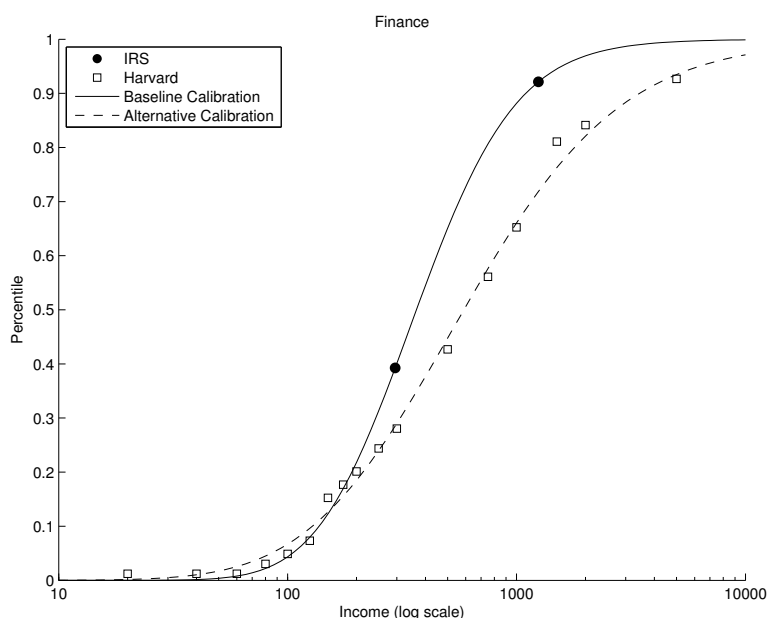


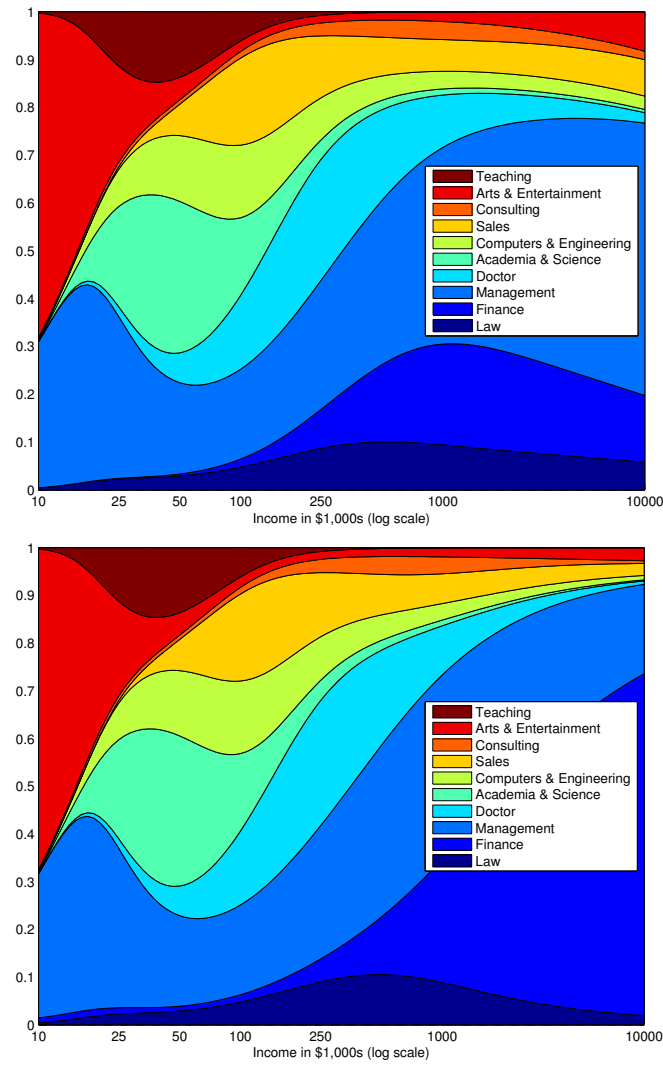
Figure 3.2: *Baseline and alternative income calibrations in Finance*



results under an alternative calibration for Finance, based solely on the Harvard data; these results are presented in full in Appendix C.1. The fit for both the baseline and alternative calibrations in Finance is shown in Figure 3.2.

Figure 3.3 shows our estimated distributions of talented individuals across professions, which is generated by normalizing the number of individuals at each income in each profession by the total number of talented individuals at that income, for the baseline and alternative calibrations. This distribution exhibits some salient and fairly intuitive features. At low incomes, most talented individuals are in Arts/Entertainment. This is resonant with the popular story of the “starving artist”. In, and only in, the lower middle class does Teaching has a significant representation. In the heart of middle income (\$25k to \$75k) the largest groups are Academia/Science and Computes/Engineering. In the upper middle class (\$75k to \$200k) Sales has a strong representation. Law and Doctor have significant representation among the the modestly wealthy (those earning roughly between \$200k and \$500k). Most of the highest income groups (those earning above \$500k) are in Arts/Entertainment, Sales or, especially, Finance and Management.

Figure 3.3: *The allocation of talent condition on income level*



Notes: The representation of talented individuals in skilled professions conditional on different income levels under the baseline calibration (left) and the alternative Finance calibration (right).

Table 3.1: Sources of externality estimates from the economics literature

	Primary Source	Method
Law	Murphy <i>et al.</i> (1991)	Cross-country regression of GDP on lawyers per capita
Finance	French (2008)	Aggregate fees for active vs. passive investing
Management	Gabaix and Landier (2008) (pro mgmt); Piketty <i>et al.</i> (2014) (anti mgmt)	Calibrated model indicating CEO pay captures managerial skill and firm characteristics (pro mgmt); Cross-country evidence that CEO pay is lower and covaries more with firm performance when taxes are higher (anti mgmt)
Doctor	—	—
Academia/ Science	Murphy and Topel (2006)	Willingness-to-pay for longevity gains of medical research
Computers/ Engineering	Murphy <i>et al.</i> (1991)	Cross-country regression of GDP on engineers per capita
Sales	—	—
Consulting	Bloom <i>et al.</i> (2013)	Randomized experiment measuring effect of consultants on plant productivity
Arts/ Entertainment	—	—
Teaching	Chetty <i>et al.</i> (2013a,b)	Future earnings of students of higher value-added teachers

3.3.2 Externality shares

To calibrate the externalities of each of these professions we drew on the fairly limited literature that tries to estimate economy-wide spill-overs from various sectors separately in different areas of the economy. We now briefly discuss our calibrations for each profession; details for selected calculations appear in Appendix C.2 and they are summarized in Table 3.1.

- Law: The only study we found of externalities from law was a cross-sectional ordinary-least-squares regression by Murphy *et al.* (1991). They investigate the impact of the allocation of talent on GDP growth rates rather on GDP levels. To be conservative and fit within our static framework, we interpreted these as one-time effects on the level

of output rather than impacts on growth rates. Attributing the effect they estimate to externalities from the talented lawyers in our sample yields $-.21$ as an externality share.

- Finance: French (2008) estimates the cost of resources expended to “beat the market” and Bai *et al.* (2013) argue that the dramatic increase in such expenditures has not made markets more informationally efficient. Viewing all of this waste as a negative externality of our talented individuals yields an externality share of $-.6$. A very similar estimate is obtained by assuming that all of the increase in the financial sector’s share of GDP (Philippon, 2013) from its trough mid-century is waste.
- Management: There is sharp disagreement in the literature over the externalities generated by managers. Gabaix and Landier (2008) and Edmans and Gabaix (2009) argue that CEO compensation is largely efficient and few externalities exist as a result, while Bertrand and Mullainathan (2001) and Malmendier and Tate (2009) argue for significant overcompensation of top managers. Following the latter, Piketty *et al.* (2014) find a $-.6$ externality share for CEOs. Because of this sharp disagreement, we consider calibrations with both $-.6$ and 0 for the externality share of Management. We denote these two calibrations by “anti mgmt” and “pro mgmt” respectively.
- Doctor: We could find no literature estimating the externality share of (non-research) medicine and so set the externality to 0 to be conservative.
- Academia/Science: Murphy and Topel (2006) estimate that medical innovation alone has generated a staggering $\$3.2$ trillion gain in welfare *each year* from 1970-2000. Even if this were the only benefit of academic research, it would translate into an externality share of 48 . Any direct use of such a number generates such a large positive externality share that it swamps everything else in our analysis. By contrast Jaffe (1989) more narrowly measures the spillover of academic research onto firm profits through excess patents in firms near universities, which leads to a much smaller 2.6 value. We settled on an intermediate value of 5 to be conservative.

Table 3.2: *Externality profiles in each of four calibrations*

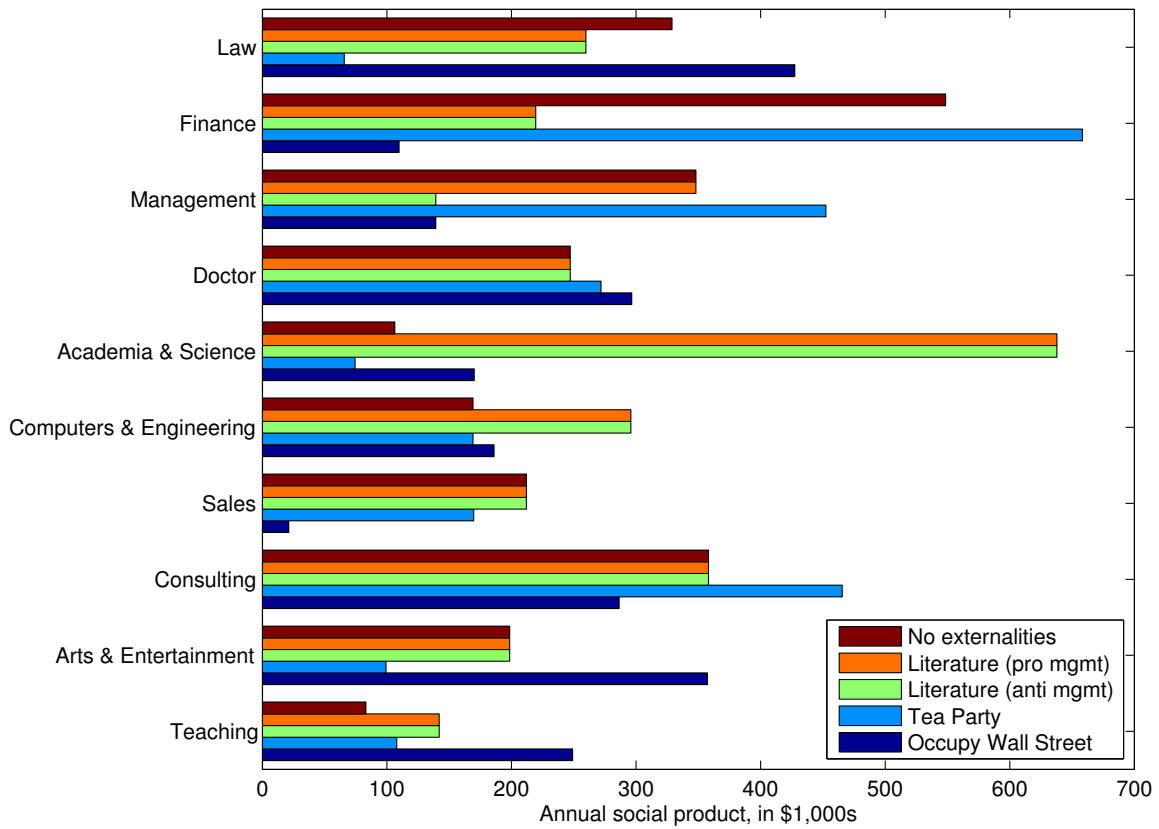
	Literature (pro mgmt)	Literature (anti mgmt)	Tea Party	Occupy Wall Street
Law	-.21	-.21	-.8	.3
Finance	-.6	-.6	.2	-.8
Management	0	-.6	.3	-.6
Doctor	0	0	.1	.2
Academia/ Science	5	5	-.3	.6
Computers/ Engineering	.75	.75	0	.1
Sales	0	0	-.2	-.9
Consulting	0	0	.3	-.2
Arts/ Entertainment	0	0	-.5	.8
Teaching	.71	.71	.3	2

- Computers/Engineering: Here again we used the analysis of Murphy *et al.* (1991), using the same methodology as in Law, and got a positive spill-over of .75.
- Sales: There is an extensive theoretical literature arguing that the welfare effects of advertising can be positive or negative, depending on whether the advertising is informative or persuasive in nature (Bagwell, 2007). The former theories imply that advertising will tend to be under-supplied in most cases (Becker and Murphy, 1993), while the latter theories suggest it will be over-supplied (Dixit and Norman, 1978). But while there have been empirical efforts to quantify the welfare effects of advertising in particular markets, such as pharmaceuticals (Rizzo, 1999) and subprime mortgages (Gurun *et al.*, 2013), we are not aware of any work attempting a comprehensive, industry-wide estimate of spill-overs, and therefore as with Doctor we assessed the externality share to be 0.
- Consulting: Bloom *et al.* (2013) conducted a field experiment to determine the causal impact of management consulting on profits. They interpreted their results as consistent with the view that consultants earn approximately their marginal product and thus we assume no externality for Consulting.

- Arts/Entertainment: While there is some evidence, and a number of good theoretical arguments, that there are some positive spillovers from the arts, we were unable to find any plausible basis for estimating the magnitude of these spillovers and they had little impact on our results (they tended to make optimal taxes slightly higher and more progressive). As a consequence, we assumed 0 to be conservative.
- Teaching: Chetty *et al.* (2013a,b) estimate the value to lifetime productivity brought by teachers of varying quality. Clark *et al.* (2013) discuss the performance of “talented” students as teachers, mostly through fellowship programs like Teach for America. Assuming that average teachers are paid their marginal product and that all additional value of high-quality teachers come through their impact on lifetime earnings, we can calculate the marginal product of an average talented teacher and divide this by their average earnings. This yields an externality share of .71.

In addition to our two literature-derived profiles of externality shares, we also want to use our model to examine positively how views about the value of various professions’ externalities can explain different views about taxation. To do so we followed two approaches. First, through discussions with with an expert on the Tea Party, Vanessa Williamson (coauthor of *The Tea Party and the Remaking of Republican Conservatism*) and several members of the Occupy Wall Street movement, we constructed externality profiles intended to represent the views of each group. In the first case the numbers assigned reflect the enthusiasm in this movement for private enterprise, skepticism of the contribution of cultural and intellectual elites and hostility to lawyers. The second position represents fierce hostility to finance and other aspects of private markets typically denigrated by the left, sympathy for legal and cultural elites’ value and enthusiasm for education. All of these calibrations are represented in Table 3.2 and in Figure 3.4, which shows what the per capita social product of various industries (from talented individuals) are under different calibrations. The latter provides a sense for the implications of these profiles for which sectors of the economy contribute most to aggregate output.

Figure 3.4: *Social Product in Different Professions*



Notes: Per capita annual social product (private product times one plus externality share) from a skilled worker in each profession, under different externality share assumptions.

Our second approach was to make an applet available online at taxapplet.appspot.com that allows readers to input their own views and receive a calibrated optimal tax policy for these views. Once sufficient data have accumulated we will attempt to analyze these patterns and their implications.

3.3.3 Results

Employing the results from Section 3.2.2, the optimal marginal tax rate at income y under ATEM is

$$(T^*)'(y) = -\bar{e}(y) - ye'(y). \quad (3.5)$$

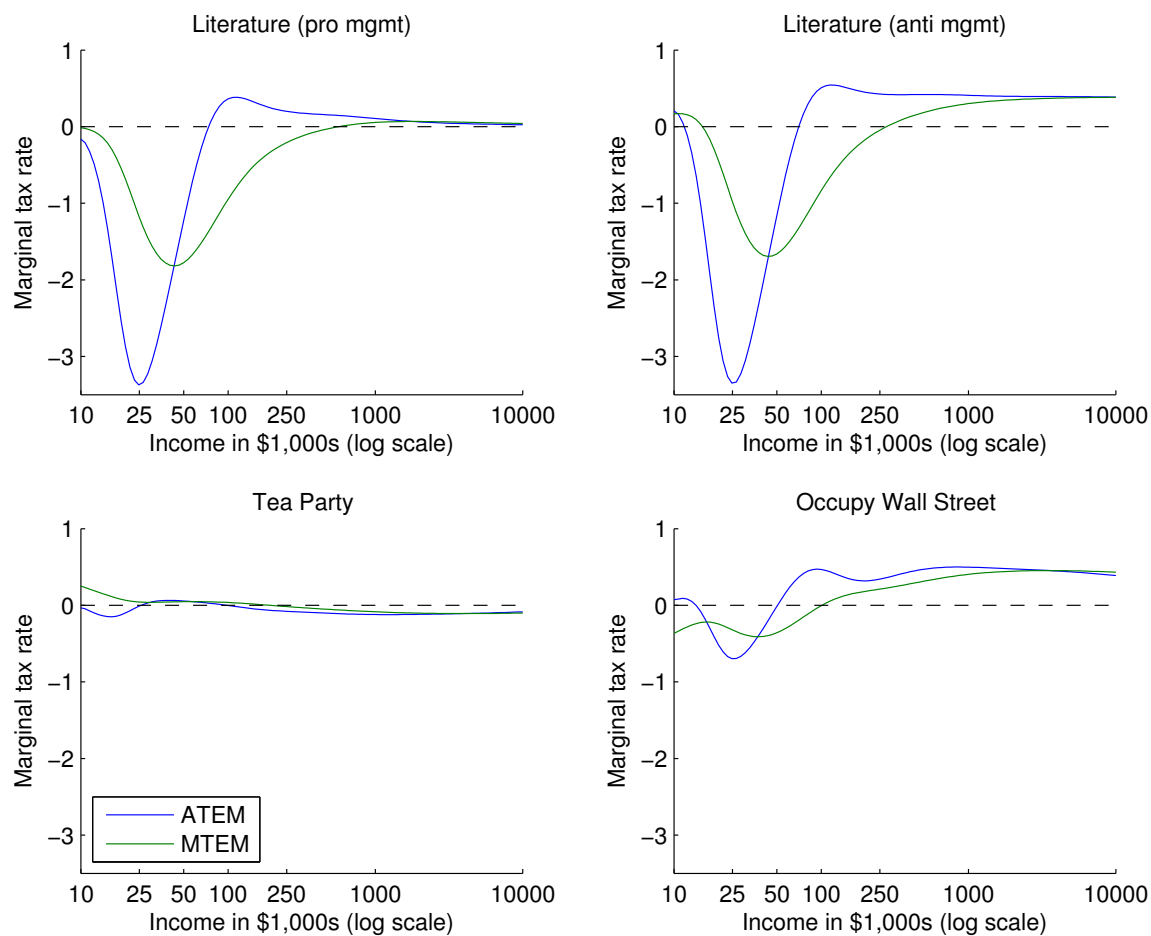
and under MTEM is

$$(T^*)'(y) = -\frac{\sum_p s_p f_p(y) e_p}{\sum_p s_p f_p(y)} \equiv -\bar{e}(y). \quad (3.6)$$

To correctly compute MTEM would require elasticities of labor supply at each income because rises in marginal tax rates will compress downwards and falls in marginal tax rates will stretch upwards the incomes earned by individuals with various average externalities. While this would not change rate schedules ordinally, it would compress the MTEM schedule across the income distribution. In the interest of comparability with ATEM, here we present the limiting case as the intensive labor supply elasticity becomes small. In this case MTEM can be computed analogously to ATEM by simply calculating the average externality at each income level and matching this to the marginal tax rate. This may also be interpreted as the direction rates should move locally beginning at current the income tax structure that induces the current by-profession income distribution. We therefore focus on this limiting case here and defer allowing a more serious consideration of labor supply elasticity until the structural model in Section 3.5.

The ATEM and MTEM policies calculated in this manner under the externality shares in Table 3.2 are shown in Figure 3.5. The first thing to note is that ATEM yields qualitatively and often quantitatively similar results to MTEM except that 1) ATEM tends to be more extreme in both directions than is MTEM for the reasons discussed above and 2) ATEM “leads” MTEM in the sense that if MTEM rates are rising (falling) in income ATEM rises first.

Figure 3.5: *ATEM and MTEM marginal tax rates*



Notes: ATEM and MTEM marginal tax rates as a function of income plotted on a logarithmic scale for our four externality share profile calibrations.

Recall that under our assumptions MTEM is also (up to an additive constant) equal to the average tax rates under ATEM. The other properties depend critically on the assumptions about externality shares employed, and thus the ATEM and MTEM policies, while looking similar given an externality profile, look radically different under the various profiles discussed above.

We begin by discussing the externality profiles from the literature. The resulting tax schedules are effectively identical up until upper middle incomes of about \$100k. Until that point, both call for massive subsidies on reaching the middle class, as the middle class hosts very positive externality professions such as Academia/Science, Computers/Engineering and Teaching while neutral or negative professions (Arts/Entertainment and Management) predominate below this. To undo these subsidies marginal rates are significantly positive (up to 50-60%) in the upper middle income range (\$50-150k). Above upper middle incomes the two schedules diverge. If Management has significant negative externalities as in the anti mgmt profile on the top right, these marginal rates of about 40% persist at high incomes, as in the Occupy Wall Street profile. If Management has no externalities, as in the pro mgmt profile on the top left, marginal rates decline to 0-5% at high incomes. Under the alternative Finance calibration, in which Finance accounts for a much larger share of top earners, marginal tax rates are close to 40% at the top even under the pro mgmt externality profile (see Appendix C.1).

Now we consider the profiles that match public views as they have intuitive (perhaps even predictable) consequences for optimal policy. First, consider the Tea Party profile. Marginal rates are slightly negative under both MTEM and ATEM after about \$200k. ATEM marginal rates dip negative, to a maximum of about -.1, between \$150k and \$300k, but are otherwise very close to zero. This is driven by the fact that we assume Tea Partiers perceive positive externalities from business and finance which account for most of upper incomes and that their presence at lower incomes is offset by the presence of academics and lawyers.

Second, consider the Occupy Wall Street profile. Under ATEM marginal rates are fairly negative (ranging from 0 to -.6) for lower and middle class individuals. They then become

and stay quite high, peaking around .5, for upper-middle and upper income individuals earning above approximately \$75k. MTEM is a less dramatic version of a similar story. Rates are modestly negative until roughly \$75k and then rise at a more moderate pace to level out around .45 for incomes in the top .1%. This is driven overwhelmingly by the negative externalities we assume Occupiers perceive managers and financiers.

3.4 Discussion

For clarity and brevity above we developed our theory with limited reference to external concepts. However, our theory is closely related to, and has implications for, several other literatures that we now discuss.

3.4.1 Allocation of talent

Baumol (1990) and Murphy *et al.* (1991) emphasize the importance of the allocation of talent for long-run growth and provide arguments and data that are the basis of some of our calibrations of externality shares. While a considerable literature builds on this analysis, little work has discussed policy tools that might be used to improve the allocation of talent. For example, Acemoglu (1995) discusses long-term cultural factors and the possibility of multiple equilibria and path dependency, but does not highlight policy tools that may be used to shift equilibria.

To get a rough sense of the importance of the allocation of talent that our taxes aim to address, consider the Goldin *et al.* data. Using our pro mgmt externality shares, we computed the (2005) income-weighted average externality share in the subset of the Goldin *et al.* data (about 80-85% depending on cohorts) who pursued careers in one of our professions.⁸ This average externality share fell from .31 to -.07 between the 1969-

⁸We only have average income data for this calculation in all professions for the later cohort. For the earlier cohort we have only Academia/Science, Management, Finance, Doctor and Law. For all of these fields except for Finance (where incomes were almost twice as high for the later cohort) incomes were very similar between the two cohorts, which is unsurprising given that they were measured in the same year. As a result, for the other professions that we did not have earlier-cohort earnings data we assumed the earlier cohort had the same earnings as the later cohort.

1972 cohorts and the 1989-92 cohorts. Suppose that this shift took place for all of our talented individuals, who capture about 43% of income.⁹ This shift would then imply a $.38 \cdot .43 \cdot 100 \approx 16.3\%$ of GDP shift in national income in some combination of reducing aggregate production and increasing the private returns of those at the top. If all were the latter (private returns rose and aggregate production stayed constant) this would account for all of the increase in the share of national income earned by the top 5% between its trough in the late 1960s and 1970s and its peak in 2008. If all were reduction in social product, it would account for around a .69% reduction in GDP growth each year for 25 years, nearly half of the 1.5 percentage point approximate reduction in US GDP growth between between the period 1948-1973 and 1982-2007. Thus, even adopting the conservative stance of ignoring the continuing shift of talent during the period 1992-2005, it is plausible that a large fraction of increased inequality and decreased growth is attributable to shifts in the allocation of talent. Of course not all of the reduction in aggregate product would necessarily have shown up in GDP, as much was a reduction in academic research and innovation (such as medical innovation) that is poorly-counted in GDP figures, as Murphy and Topel (2006) highlight. Nonetheless, it seems plausible that the changes in the allocation of talent during this period alone had large aggregate and distributive consequences for welfare, on the same order of magnitude as the widely discussed shifts observed in aggregate statistics.

3.4.2 Labor supply elasticity debate

An important recent controversy in public finance concerns the elasticity of labor supply. One literature, surveyed by Saez *et al.* (2012), highlights that the short-term elasticities of labor supply and taxable income that can be measured through natural experiments are low (on the order of .1 – .5 for high income earners) and largely driven by welfare-irrelevant evasion or inter-temporal substitution. A second literature, surveyed by Keane (2011) and

⁹It seems likely that that the shift among Harvard alumni was more extreme than in the rest of the population. However, it also seems likely that the shift among these students was greater by the mid-2000's than it was by 1989-92, and greater from the mid-1950's than it was since the late 60's. Here we assume that these two effects balance.

Keane and Rogerson (2012), argues that this evidence is consistent with large long-term elasticities of labor supply that would not appear in short-term estimates and that such large long-term elasticities also help rationalize international data. The first literature concludes that the deadweight loss from taxation is small due to low labor supply elasticities and argues for highly redistributive tax policy; the second literature argues that taxes are highly distortionary.

Our work offers a third view. On the one hand, our calibration of career-switching elasticities is largely consistent with Keane and Rogerson arguments for high (perhaps even higher than they claim) long-term labor-supply elasticities among talented and wealthy individuals. On the other, it suggests that high long-term elasticities (interpreted as career-switching elasticities) need not imply lower optimal marginal tax rates. In fact, they may make optimal marginal rates higher over many ranges if high marginal rates are implemented for efficiency rather than equity reasons.

Thus it emphasizes that the externalities of various professions, rather than the elasticities of labor supply, may be the primary determinants of optimal income taxation. If externalities are small or more positive at higher incomes, given our calibrated career-switching elasticities, optimal tax rates are likely to be low even with a standard redistributive motive as we discuss further in Subsection 3.5.3. On the other hand, if externalities are large and more negative at higher incomes, optimal tax rates on the wealthy are likely to be very high, increasing in the long-term elasticity and not very sensitive to redistributive motives.

The relative empirical attention these two sets of parameters have received seems disproportionately-weighted towards labor supply elasticities. Thousands of papers, a few hundred of which are surveyed in the papers discussed above, have been devoted to measuring the elasticity of labor supply and of taxable income more broadly. On the other hand, the few papers we use to generate our literature-based externality profiles represent, to our knowledge, the extent of research on industry-specific externalities. Our theory suggests this topic deserves more attention in future research.

3.4.3 Debates on taxation outside neoclassical economics

Public debate over tax policy rarely focuses on the parameters emphasized in optimal tax theory (viz. the degree of inequality or the responsiveness of work to taxes). Instead, as Mankiw (2010) notes, much rhetoric instead focuses on whether the rich “deserve” the wealth they have accumulated. The left attempts to delegitimize the wealth of the rich (claiming it is misbegotten or crooked) and of the right to hold the wealthy up as job creators and entrepreneurs. Some simple public opinion data (Parker, 2012) is suggestive here. While 55% of Republicans believed the rich were more likely than others to be hardworking and 18% believed they were more likely to be honest, only 33% and 8% of Democrats agreed on each count respectively. On the other hand only 42% of Republicans believe the rich are more likely than others to be greedy, while 65% of Democrats do. According to a Quinnipiac University poll (Brown, 2011), Republicans believe 59% to 28% that public sector workers are overpaid while Democrats believe 38% to 31% that they are underpaid. A CNN poll CNN/ORC (2011) found that while an equal number of Tea Party supporters had “a great deal” or “some” confidence that Wall Street Bankers acted in the interests of the overall economy as had no confidence at all, 73% of opponents of the Tea Party had no confidence while only 13% had some or a great deal of confidence. Similarly while Tea Party supporters only believe 58% to 40% that these bankers are overpaid and believe 50% to 47% that the bankers are not dishonest, Tea Party opponents believe 89% to 10% that they are overpaid and 75% to 23% that they are dishonest.

These opinion patterns are not new: in the nineteenth century reformers accused the wealthy of being “robber barons” and Marx (1867) accused the wealthy not of being unresponsive to tax rates but of being unproductive exploiters of the truly productive working class. Literature on the political right, such as Rand (1957), provide hagiographic representations of the social contributions of the wealthy rather than depictions of their willingness to shirk if taxes rise. In fact that book is largely devoted to the unwillingness of the rich to shirk even when they are nearly enslaved, a setting that would under the Mirrlees theory be ripe for a high tax rate.

Our theory provides a natural, formal, quantitative language for these debates which fit only unnaturally with the Mirrlees framework. In so doing, data can be brought more easily to bear; disputes become questions of which professions exactly are claimed by different sides to have different degrees of externalities and which professions in fact earn which incomes. Such questions should be easier to settle empirically than are broad and vague claims about the moral worth of different social groups. At the very least it suggests that if economics wants to speak to these public debates it should focus more attention on the extent and degree of these externalities.

3.4.4 Closely related literature

The work most closely related to ours is that of Philippon (2010), Rothschild and Scheuer (2014a) and Rothschild and Scheuer (2014b).¹⁰ All three papers investigate public policies aimed at reallocating talent. Philippon considers the use of taxation to affect the allocation of talent between a financial and entrepreneurial sector in a model where financiers serve as conduits for funding of research and entrepreneurship. By contrast, externalities are directly imputed (rather than arising endogenously) in our framework and we consider many professions simultaneously. Furthermore Philippon focuses on sector-specific instruments rather than horizontally equitable policies.

Rothschild and Scheuer are closer to our work in this dimension, considering horizontally equitable optimal tax policy. Our papers focus on different aspects and approaches to this problem. In Rothschild and Scheuer (2014a) they restrict attention to two professions (one with negative externalities and one with none), but in both papers allow for richer targeting of externalities: in their model, rather than externalities accruing uniformly across the population they may be targeted either at individuals within the “rent-seeking” sector or towards individuals in the productive sector. In that paper, like Philippon, they emphasize the perhaps counter-intuitive theoretical finding that it may be theoretically optimal for

¹⁰The original version of the paper that this draft was based on, “Psychic Income, Taxation and the Allocation of Talent”, was written prior to the publication of the first paper and a draft of the second. However we have no reason to believe the authors of either paper were aware of that work.

policy to (implicitly or explicitly) subsidize the unproductive wealthy, in their case because the negative externalities within the unproductive sector discourages further wasteful entry into this sector. In Rothschild and Scheuer (2014b), they consider a far more general case that nests our analysis as a special case and provide general formulae for optimal taxation, as well as qualitative, directional results and uncalibrated empirical sufficient statistics in some special cases that do not nest our model, such as a model with only one profession generating an externality.

By contrast our emphasis is quantitative: our goal is to determine how the magnitude of optimal taxes varies with assumptions about externalities, allowing for the many professions that make up real-world top income distributions and for positive as well as negative externalities.¹¹ Given the lack of empirical information we have on the targeting of externalities, this requires us, among other restrictions, to either ignore the effects highlighted by Philippon and Rothschild and Scheuer or to treat targeting as another subjective input. We chose the former course because we believe that assuming externalities are born evenly across the economy is a reasonable benchmark. For example, while it is true that arbitrage activity may reduce the returns to some financiers, more high speed traders require other investors to make greater investments to avoid being front-run, thereby hiring more financiers. The legal profession hosts similar arms races. Most of the positive externalities of academic work and teaching are born broadly by the public and not by any particular sector. Furthermore, the dramatic rise of the financial sector in the United States over the last thirty years as documented by Philippon and Reshef (2012), Goldin *et al.* (2013) and Philippon (2013) seems to belie the equilibrating forces Philippon (2010), Rothschild and Scheuer emphasize. Thus we adopt what Rothschild and Scheuer call the “naïvely Pigouvian” perspective that wages, externalities and entry into professions are independent. An interesting extension of our work and theirs would be to incorporate the targeting emphasized by their work into a

¹¹As a result of these contrasting goals, we take a very different methodological approach to Rothschild and Scheuer as well. For example, we use the taxation principle in the spirit of Saez (2001) rather than the revelation principle of Mirrlees (1971). Similarly we make assumptions (about substitution patterns across professions) that are convenient for calibration to the data we have available, rather than for illustrating possibilities.

quantitative, calibrated framework such as ours.

3.5 Structural Model with General Ability

In Section 3.2 we developed an intuitive characterization of optimal implicit Pigouvian taxation of income under a stark set of assumptions. In this section we discuss the robustness of the results to adding realistic features to the model using a combination of analytical and computational techniques. While our characterization is less clean in these cases, our qualitative results are actually strengthened, and we are able to evaluate our analysis quantitatively in greater depth.

The model in Subsections 3.2.2 and 3.3.3 assumes a strong form of orthogonality between elasticities and externalities conditional on income level. One natural scenario in which this would be violated (along the career-switching margin) is if each individual possesses a general level of ability which determines her wage in all professions. In this case, individuals systematically shift from professions with high average earnings to those with low average earnings when taxes rise. For example, an academic making a million dollars a year—near the top of the distribution in that industry—would earn far more in finance; she is just exceptionally talented. A marginal tax rate increase at any point between these two salaries would only reduce the relative attractiveness of finance (and a marginal tax rate increase outside of this range would not change the relative attractiveness of the professions at all). On the other hand, a financier earning a million dollars a year—much lower in the finance salary distribution—is not so exceptionally talented, and would face a much lower income in academia. Therefore a marginal tax rate increase between between these two salary levels would render academia relatively more attractive—potentially inducing the financier to become a professor. Intuitively, both individuals face the same *ranking* of wages in the two professions, and thus marginal tax rates are likely to cause only switches down from generally high-paying to generally low-paying professions, not the reverse.¹²

¹²The same scenario is also likely to imply differential intensive elasticities of hours supply *conditional on income*. An academic making a million dollars a year is likely to be working many more effective hours than is

While few simple analytical results may be obtained in general when we depart from this story, under a stark simplification of the general ability model a natural result is possible. In particular, suppose that ability is fully portable across careers in the sense that individuals switch between the same quantile of an underlying “reference distribution” of profession-conditional income when they switch, that there are exactly two professions and that the profession-contingent income in these two are strictly ranked by first-order stochastic dominance. Then marginal tax rates will always be more sensitive to externalities than those given by ATEM in the sense that, starting at ATEM tax rates, there is a first-order welfare gain from raising marginal tax rates at every point if the high-paying profession has more negative externalities and a first-order gain from lowering them if the high-paying profession has more positive externalities. We state and prove this result formally in Appendix C.3 and focus here on a computational structural model of many professions incorporating this feature.

3.5.1 Calibration

To formalize the general ability idea, suppose that every individual i is endowed with a uniformly distributed general ability $a_i \in [0, 1]$. She also receives a vector of non-pecuniary total costs of work in each profession ϕ_i forming a type $\theta_i = (a_i, \phi_i)$. Individual i earns wage $G_p^{-1}(a_i)$ in profession p , where G_p is a “reference distribution” of wages in profession p . Additionally, in order to calibrate a full structural model we

1. Assume away or impose a uniform value from the literature for intensive elasticities,
2. Impose the common logistic functional form on non-pecuniary utilities,

an financier earning that salary both because her wage in finance is likely to be lower (and thus, to be earning the same income, she must be working harder) and because her non-pecuniary cost of work is likely to be lower. She is therefore likely to be pressing herself to the limit of her exertions and thus to be much less elastic to an increase in her monetary compensation than is the financier. Thus one should expect individuals in a high-paying profession to be more elastic, *conditional on income*, to compensation than are those earning the same income in a generally low-paying profession. In a previous draft of the paper we had a formal result to this effect. However, because our focus is not on the intensive margin we have omitted this result.

3. Calibrate the free parameter of this distribution to the career-switching elasticity estimated in Section 3.3 above,
4. Then jointly estimate the mean utilities of different professions and a non-parametric reference income distribution to match the shares choosing professions and the profession specific income distributions as estimated in Section 3.3.1.

We then show, as suggested by the results above, that the results of our simple model in Section 3.3.3 are generally strengthened in this more realistic setting. We now discuss these steps; they are described in technical detail in Appendices C.3 and C.4.

Each individual is endowed with a fully-portable, general ability $a_i \in (0, 1)$ that would entitle her to pecuniary income from that quantile of the reference income distribution in each profession. Each profession has a mean non-pecuniary income and each individual has, in addition to this an independently and identically distributed (across professions and individuals) Type I Extreme Value component of her non-pecuniary income in each profession. The scale parameter of these idiosyncratic draws, $\beta(a_i)$, determines how responsive individuals are to changes in the attractiveness of different professions. $\beta(a_i)$ is assumed to be proportional to the mean income that an individual of ability a_i would earn under, before taxes, across all professions she might choose. This assumption loosely reflects the fact that individuals who earn higher incomes are likely to put relatively more (dollar) value on the non-pecuniary benefits of their professional choice, as otherwise we would systematically conclude that high-ability individuals uniformly sort into the highest paying professions while lower-ability individuals make more diverse choices. In any case we have checked that our results are not quantitatively significantly sensitive to several alternative specifications of $\beta(a_i)$, including it being constant.

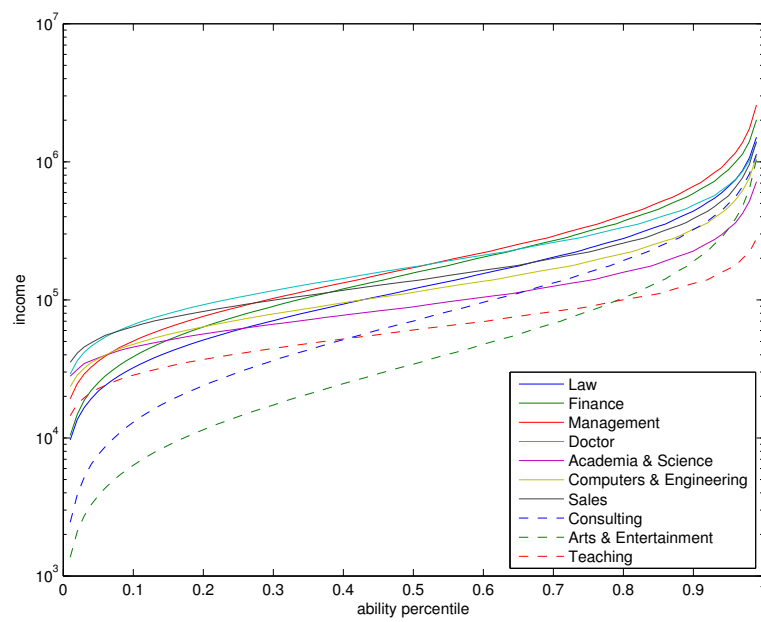
To calibrate the β , we need a measure of mobility across professions in response to relative income shocks. Intuitively, if β is large, allocation across professions is primarily driven by individual idiosyncrasies and thus there will be little individual response to changes in the relative wages of different professions and conversely if β is small. To estimate β , we use the elasticity of individuals entering the finance profession in response

to changes in relative wages discussed at the end of Subsection 3.2.2 above; see Appendix C.4 for details. As a baseline we use the assumption that the migration of individuals into finance between the 69-72 and 89-92 Harvard cohorts respectively was driven by a \$150k change in relative wages in 2005 dollars. This is sufficient to calculate the value of β as it is a scalar. We then check robustness to using an increase of either \$75k or \$300k in Appendix C.6. We show that our qualitative results are robust to this assumption, but the larger the elasticity (i.e., the lower the assumed income change that induced the migration), the larger are the welfare impacts and the importance of externalities.

Finally, we jointly recover the mean psychic income and reference income distributions so as to exactly rationalize a discrete grid of the observed distribution of individuals across professions and of income conditional on professions estimated in Subsection 3.3.1. We do this using a numerical non-linear equation solver as described in Appendix C.4.

The resulting reference income distributions are shown in Figure 3.6. These indicate the salary that an individual from a given quantile of the ability distribution would earn in each profession, and correspond closely to intuition. Above median ability, Finance and Management are essentially tied for the most lucrative career. Below median ability Doctors, and at the very bottom Sales, do better. This corresponds to the common intuition that being a doctor or “mad man” is a less risky career path than finance and management are. Arts/Entertainment is the worst paying profession except at the top, where it is quite lucrative, exhibiting the well-known superstar structure of that profession (Rosen, 1981). Teaching and, to a lesser extent, Academia are also quite low-paying, while Law and Sales are on the more lucrative side, with Law having greater inequality than Sales. Interestingly Consulting is the most unequal profession, scoring near the top among the very able but second to last at the bottom of the income distribution. Given that these results are qualitatively similar to the equilibrium, observed income distributions of Subsection 3.3.1 to which they are matched, we found in a previous draft very similar results if the empirical distributions are simply used directly.

Figure 3.6: *Reference Income Distributions*



Notes: Reference income distributions by percentile in the various professions, the income that an individual in that quantile of the ability distribution would earn if she worked in the relevant profession.

We also estimate the income distribution of the US working population outside of the talented individuals considered above. We denote this group “unskilled”. Because we will assume that the externality of this group is 0, its income distribution was not necessary to compute the optimal tax schedule given by equations 3.5 and 3.6. We compute its population share as the residual of the other shares, and we calibrate its Pareto-lognormal parameters so that the total income distribution (across all professions) matches three moments: the mean population income, and the 99th and 99.9th income quantiles.¹³

3.5.2 Optimal tax rates

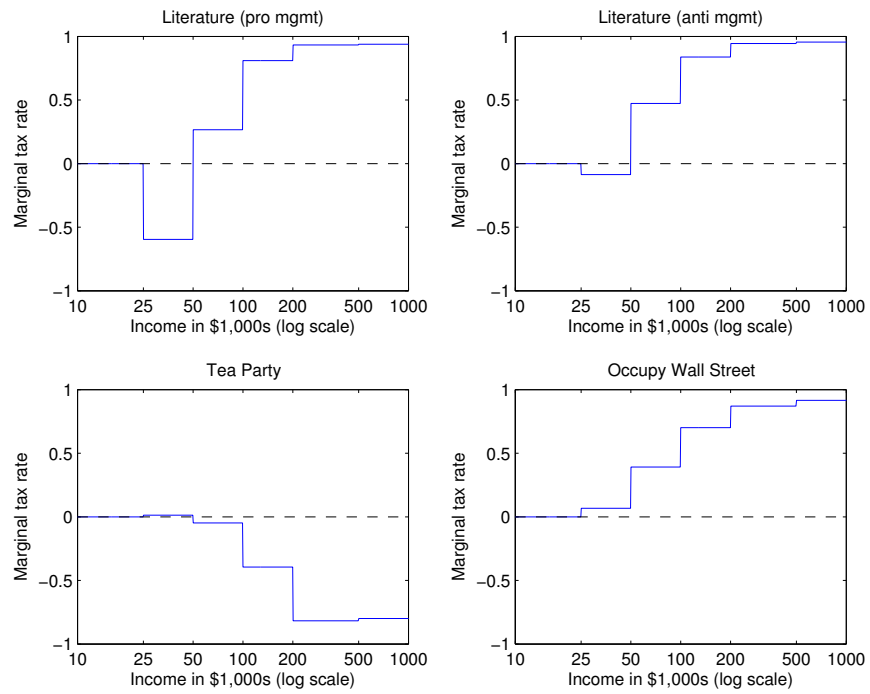
This model obviously has many limitations: ability is fully portable, disallowing comparative advantage which empirical evidence suggests is clearly important (Kirkebøen *et al.*, 2014); substitution patterns across professions obey the implausible independence of irrelevant alternatives assumption; there are no intensive elasticities, etc. Perhaps worst of all, exactly how these assumptions impact our results is far from transparent. A more realistic model of any of these features would require much more detailed data than we have access to. Here our goal is to provide an example model with more than two professions where under the general ability assumption we obtain a result in the spirit of that demonstrated in general in the simpler case: that optimal marginal rates tend to be more responsive to externality profile assumptions than in our baseline ATEM analysis.

To do this, we searched computationally for optimal piecewise-constant marginal rates with a zero marginal rate applying on income from \$0 to \$25k and optimal rates estimated for \$25k-\$50k, \$50k-\$100k, \$100k-200k, \$200k-500k and \$500k and above brackets.¹⁴ We focused on these piecewise structures both because they are realistic and because they

¹³In principle, we could non-parametrically estimate the “unskilled” income distribution by subtracting the estimated CDFs of the skilled occupations from the the IRS empirical distribution of income. The problem is that “unskilled” is not constrained to look like a Pareto log-normal at the while the other occupations are. As a result, “unskilled” picks up all the mass at the (above \$10,000,000), which seems implausible. Thus we also estimate “unskilled” as a Pareto log-normal.

¹⁴Given the minuscule share of skilled agents earning less than \$25k, the marginal tax rate between \$0 and \$25k is approximately collinear with the budget balancing demogrant and generating a multiplicity of optima; we avoid this issue by constraining the lowest marginal tax rate to zero.

Figure 3.7: *Optimal Tax Rates in Structural Model*



Notes: Computational optimal piece-wise-constant marginal tax rates for a six-bracket tax system ($[0, 25k)$, $[25k, 50k)$, $[50k, 100k)$, $[100k, 200k)$, $[200k, 500k)$ and $[500k, \infty)$) in our computational general ability model for the listed externality profiles.

significantly reduce the computational costs of searching for an optimum compared with a fully non-linear setting given the absence of a closed-form expression for the optimum. To avoid incentives for money burning we introduced a very small, uniform intensive labor-supply elasticity of .01. Instead simply bounding rates between -100% and 100% yield similar results, but ones that commonly hit the boundary.

These are pictured in Figure 3.7. The results are qualitatively similar to those we obtain in Subsection 3.3.3, but with smaller subsidies to the middle class in the Literature and Occupy Wall Street profiles, and more extreme marginal tax rates at high incomes. The Tea Party profile in the bottom-left panel of Figure 3.7 as in Subsection 3.3.3 lead to negative marginal rates (except on the lower middle class), but here they are much more negative than before, reaching around -75% rather than -10% for the top bracket and being nearly -50% starting as low as \$100k.

Occupy Wall Street and both Literature profiles call for extremely progressive tax regimes. Occupy, in the bottom-right panel, calls quite high marginal tax rates beginning at \$50k at around 40% and rising to near confiscatory levels at \$100k and above.

The literature-based calibrations, in the top row are similar except being even more confiscatory for the wealthy and supporting a subsidy for the middle class. However these are modest compared to those called for under ATEM in Section 3.3.3, since such subsidies are only useful in inducing employees in lower paying professions (which tend neutral externalities, such as art) to switch to middle-class incomes. Since these potential switchers do not create substantial negative externalities, and would not be terribly productive in the positive-externality producing professions in the middle class, large subsidies are not worth even the very small distortion under a labor supply elasticity of 0.01. This provides suggestive evidence that the effects discussed analytically at the beginning of this section above may be quantitatively important and thus that our main results from Subsection 3.3.3 may significantly *understate* how widely optimal marginal rates vary as a function of externality shares, except in the case of marginal subsidies to join the middle class.

3.5.3 Quantitative importance of elasticities vs. externalities

We now seek to compare the quantitative importance of externalities relative to the labor supply elasticity, which has been the focus of the traditional optimal tax literature following Mirrlees (1971). To do this, we introduce a redistributive motive and intensive margin labor supply elasticity into our model. We select a particular baseline case: a labor supply elasticity of 0.25 and zero externalities, then we ask which of two variations around this baseline causes a greater change in the optimal tax code: varying the labor supply elasticity across a reasonable range of estimates drawn from the extensive literature estimating that parameter, or introducing our various externality profiles. Our results suggest externalities are at least as important as are labor supply elasticities to optimal marginal tax rates, at least given the current plausible range of uncertainty on these.

For these results, we begin with our computational model developed in the previous subsection. To add a redistributive motive, we follow the strategy of Diamond (1998) by placing different social marginal welfare weights on individuals while continuing to model their utility as quasi-linear in income. Also following Diamond, we generate these marginal social welfare weights based on their quasi-linear utility, assuming a constant relative risk-aversion of .95 so the marginal social welfare weights are approximately the inverse of the relevant income.¹⁵ To introduce an intensive-margin labor supply elasticity, we set ϕ to take a power form with exponent that is homogeneous across individuals so that all individuals, regardless of hours worked, have a constant labor supply elasticity as in Saez (2001).

We then run a “horserace” over the importance of labor supply elasticity as compared to externalities in determining optimal taxes.¹⁶ To do this, we first take a baseline model with log-utility-based redistributive preferences, a uniform intensive labor supply elasticity of .25 (as seems roughly consistent with the literature as surveyed by Saez *et al.* (2012)), our calibrated career-switching elasticities and no externalities. Rates below \$25k are again constrained to 0%. The results for optimal taxes in this calibration are pictured in the

¹⁵For computational reasons this is more convenient than a literally logarithmic, inverse-weights analysis.

¹⁶We are grateful to Matthew Weinzierl for suggesting this analysis.

far upper right panel of Figure 3.8. A few features of this baseline are prominent. First, marginal rates are moderately high throughout the income distribution, ranging from a bit under 35% to 65%. However, marginal rates decline at the top of the income distribution, despite the fact that our calibration shares with Saez (2001) Pareto tails on the income distribution, which usually leads to monotonically increasing marginal rates at the top. The reason is that career-switching becomes very important at high income levels and, when it occurs, generates a discrete drop in tax revenues, unlike the intensive margin adjustments considered by Saez. This is qualitatively consistent with the intuitions of Keane (2011) and Keane and Rogerson (2012) discussed in Subsection 3.4.2. Veiga and Weyl (2013) show that, in spite of Pareto tails to income distributions, if the fraction of current income individuals can earn by an outside option (such as career switching) is affiliated (Milgrom and Weber, 1982) to income then optimal tax rates are regressive.

We then consider two ways of varying parameters about this baseline. First, in the left collection of panels, we add to this model the externality views of our calibrations from Subsection 3.3.2. The combination of redistributive preferences and a moderate intensive labor supply elasticity significantly moderate the extremity of optimal tax rates across calibrations. However, comparing the two extremes, the Tea Party and Occupy Wall Street, optimal tax rates still differ quite dramatically. Rates under the Tea Party profile are moderate and decline at the very top. Marginal rates run from approximately 40% on the lower middle class up to nearly 60% on the wealthy but then fall back down to 50% on the very wealthy. By contrast, desired taxes of the Occupy movement movement are highly progressive up to the wealthy and plateau for the very wealthy. The optimal rates based on the Literature profiles are quite similar to Occupy Wall Street, except that they are even flatter at the very top under anti mgmt and decline slightly under pro mgmt. Under the alternative finance calibration, marginal rates are monotonically increasing even for pro mgmt (see Appendix C.1). The differences across profiles are even more striking in terms of average tax rates, which are not shown here.

The second way we vary about the baseline model is by changing the intensive labor

supply elasticity, between .1 and 1, a range that seems to span reasonable beliefs for the intensive margin according to both Saez *et al.* (2012) and Keane (2011). At .1, rates are higher (ranging from 50% to above 80%), but drop significantly at the top, falling from over 80% to about 70% for the very wealthy. For an intensive elasticity of 1, rates for the middle class are quite low (10%) and rise to around 40% at high incomes.

Which affects optimal tax rates more, a reasonable range of labor supply elasticities or a plausible range of externality views given current public debates? While the average rates move more with intensive labor supply elasticity than they do under the externality calibrations, the qualitative nonlinear structure of optimal taxes moves more with externalities than with labor supply elasticities. In particular, marginal tax rates decline notably at the top under Tea Party views and plateau under the Literature and Occupy Wall Street profiles, while they decline moderately for each of the intensive labor supply elasticities. Therefore, over reasonable ranges of public debate, externalities may have approximately as large an impact on the level and perhaps an even larger impact on the structure of optimal taxes as plausible views about labor supply elasticities. This reinforces our argument in Subsection 3.4.2 that these externalities merit much closer attention in empirical work relative to the large literature devoted to labor supply elasticities.

3.5.4 Quantitative welfare gains

For our final computational exercise, we analyze the welfare gains from optimal policy as compared to the *laissez-faire* regime with zero marginal tax rates everywhere under various externality calibrations and compare these to the welfare gains that could be achieved if profession-specific optimal taxation were possible. We consider this question under the model of Subsection 3.5.2 above which features no intensive elasticity or redistributive motive. In addition to reporting the impact on social welfare we also consider the impact on aggregate pecuniary income (viz. GDP), which excludes changes in psychic income.

As with all other results, our results here, which appear in Table 3.3, differ markedly depending on the externality profile used. Welfare gains relative to *laissez-faire* from the first-

Table 3.3: Welfare Gains

	Gains from 1st best in Social Welfare (GDP)	Gains from 2nd best	Share of potential from 2nd best
Literature (pro mgmt)	17% (32%)	1.5% (-.15%)	8.8% (-.45%)
Literature (anti mgmt)	21% (39%)	3.3% (1.8%)	16% (4.8%)
Tea Party	1.1% (4.3%)	.56% (3.7%)	49% (85%)
Occupy Wall Street	4.1% (5.3%)	1.7% (-.071%)	42% (-1.3%)

Notes: Quantitative welfare gain, compared to laissez-faire (zero marginal tax rates everywhere) in our baseline computational model from optimal non-linear taxation (2nd best) and optimal profession-specific taxation under various calibrations.

best, profession-specific taxation optimum are largest under the profiles from the literature as the externalities are most extreme there; they are smallest under the Tea Party's views because these have the smallest externalities. The first thing to note is that the numbers are quite large. Gains from moving to the first-best are around 20% of social welfare and 30-40% of GDP under the literature calibrations.

Furthermore, we have restricted this model so all gains come only from career switches, which occur among the less than 4% of the population we label talented. The private product of these individuals is 43% of GDP. Thus, under these views, the reallocation of these individuals across professions induced by first-best, profession-specific taxation could nearly double the output from these individuals. These gains are one to two orders of magnitude larger than those accruing from capital taxation in a dynamic optimal income taxation model according to the baseline calibration of Farhi and Werning (2012), a topic that has attracted far greater attention in the literature than have the externalities we consider.

Gains under the Occupy views are still large, but more modest; they are fairly small under the Tea Party views.

The fraction of the first-best gains from optimal, profession-specific taxation achieved by the second-best optimal horizontally equitable non-linear tax is quite small under the literature views. 9-16% of welfare gains are achieved and GDP actually falls under the pro mgmt views because individuals migrate towards more rewarding careers at the cost

of reduced output, even including externalities. Under the anti mgmt literature views, GDP increases by about 2% in the second-best, which is about 5% of the GDP gains under the first-best. These gains are a very low fraction of the first-best because there is so much heterogeneity of externalities conditional on income for these externality profiles. In Appendix C.7 we show that while under the first-best 95% of talented individuals go into Academia/Science under the second-best they increase their presence in Academia/Science only moderately (from 10.6% to 15.6%), while distributing themselves more uniformly among less lucrative careers than under *laissez-faire*. For example, the not-especially-productive field of Arts & Entertainment more than doubles in size under second-best policy. Under Tea Party and Occupy Wall Street profiles, by contrast, there is a much-closer alignment between income and externalities and thus a greater fraction (40-50%) of social welfare gains from the first-best are achieved at the second-best.

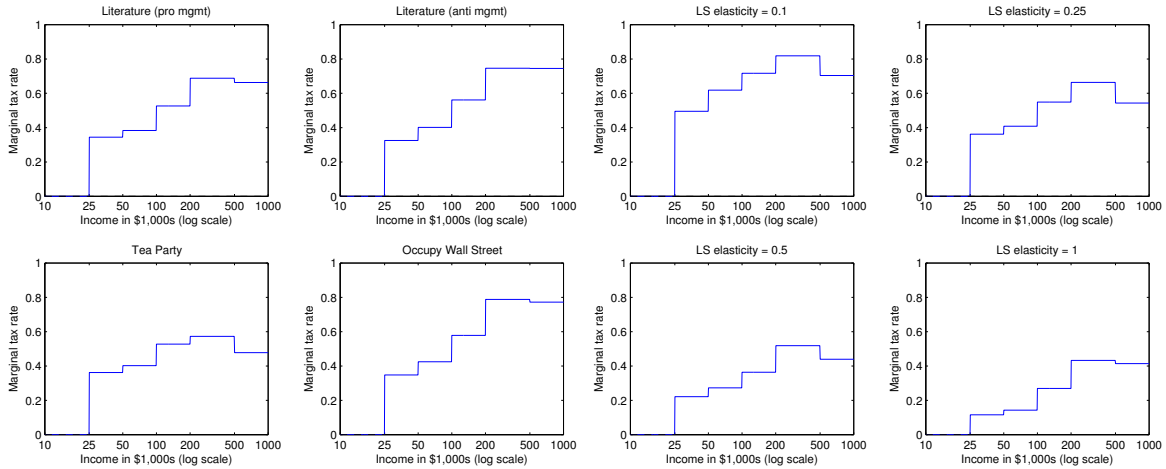
3.5.5 Effects of the Reagan tax reforms

In the previous section we considered the impact of optimal (first- or second-best) tax reforms relative to *laissez-faire*. These are abstract exercises in the sense that they do not correspond to plausible ranges of policy reforms enacted or considered in the United States. Therefore, in this section, we consider in the context of our model the impact of a policy reform that actually occurred: the Reagan tax reforms that transformed the US income tax system from 1980-1990.

To do so we use The Tax Foundation (2013)'s history of US federal income tax brackets. While these are obviously an imperfect approximation to the full wedge faced by each income, we believe they provide a reasonable sense of the change in these wedges over the 1980-1990 period. We compare results predicted by our model under these two tax regimes along four dimensions: the allocation of talent, social welfare, GDP and top income shares.

First consider the impact of the reforms on the allocation of talent, as shown in Table 3.4. In relative terms the largest impacts were to move talented individuals out of Arts/Entertainment, Teaching, Consulting and Academia/Science into Finance and Man-

Figure 3.8: *Horserace between elasticities and externalities*



Notes: A “horserace” between the importance of labor supply elasticities and externalities in determining optimal taxes. The left set of panels hold fixed redistributive motives at log-utility with an intensive labor supply elasticity of .25 and our calibrated career-switching patterns, but allow externalities to vary across our calibrations. The right set of panels hold externalities fixed at 0 and other features fixed as in the other calibrations, but allow the intensive margin of labor supply elasticity to vary from .1 to 1.

Table 3.4: *Reallocation of Talent from Reagan Tax Reforms*

	Pre-Reagan	Post-Reagan	Absolute Change	Relative Change
Law	6.6%	6.4%	-.2%	- 3.2%
Finance	5.1%	6.0%	.9%	18%
Management	22%	25%	3.0%	14%
Doctor	16%	17%	1.4%	6.3%
Academia/Science	14%	12%	-1.5%	-11%
Computers/Engineering	12%	11%	-.9%	-3.2%
Sales	13%	13%	-.2%	-1.7%
Consulting	2.4%	2.0%	-.4%	-17%
Arts/Entertainment	3.9%	3.0%	-.9%	-23%
Teaching	5.8%	4.7%	-1.1%	-19%

Notes: Estimates from our baseline calibrated structural model of the effects of the Reagan tax reforms and the allocation of talent.

agement by 10-25% of the pre-reform magnitudes in these fields. In absolute terms the largest reallocations were from Academia/Science and Teaching into Management, Doctor and Finance.

Thus in our model the Reagan tax reforms significantly shift individuals from low-paying towards high-paying occupations. Under our calibrations based on the literature, especially the anti mgmt calibration, this causes a significant reduction in welfare. Under pro mgmt welfare falls by .56% and GDP rises by .083% while under the anti mgmt calibration welfare falls by 1.3% and GDP by .77%. In either case these “supply side” tax cuts leave GDP essentially unchanged or lower it, compared to the results under the Tea Party calibration where GDP rises by 2.1% as a result of the reforms.

Perhaps most interesting is the impact of the reforms on top income shares. As Hacker and Pierson (2010) summarize, many economists have argued that, while the Reagan reforms may have increased post-tax income inequality, it is unclear how they could have been responsible for the large rises in pre-tax income inequality in the past 35 years documented by Piketty and Saez (2003). Our analysis provides a clear mechanism for this: lower taxes cause substitution between psychic and pecuniary income, leading both to an increase in the earnings of the very talented and, at least under the calibrations from the literature, a reduction in net positive spillovers to other income groups. Both factors tend to increase the relative pecuniary income of the very wealthy.

To measure the size of this effect, we compute the increase in income share accruing to the top 1% of the population in our model, before and after the Reagan reforms. We compare this to the corresponding increase in income for skilled professions in Bakija *et al.* (2012), which is broadly consistent with the patterns documented by Piketty and Saez, though muted by the exclusion of capital gains and earnings in other professions. According to Bakija *et al.*, the share of total income received by members of our skilled professions in the top 1% of the population income distribution rose 3.9 percentage points from 1979 to 1993 and 1.4 more percentage points by 1998. In our model the Reagan reforms cause a .8 percentage point rise in the same figure. Thus, according to our model, the

Reagan reforms would account for 21% of the rise in the pre-tax share of the top 1% among talented individuals from 1979 to 1993 and 15% of the rise through 1998. As with our other quantitative findings, these results are far stronger (about 2-2.5 times as large) under our alternative calibration of the share of talented individuals in Finance, based on fitting more closely the Harvard data, as discussed in Appendix C.1.

3.6 Conclusion

This paper proposes an alternative framework for the optimal taxation of top incomes to the standard redistributive theory of Vickrey (1945) and Mirrlees (1971). Income taxation acts as an implicit Pigouvian tax that is used to reallocate talented individuals from professions that cause negative externalities to those that cause positive externalities. Optimal tax rates are highly sensitive to these externalities. If they are as large as our reading of the literature suggests, the worsening allocation of talent in the United States is large enough to account for all of the increase in inequality or nearly half of the fall in growth between the 1948-1973 period and the 1982-2007 period.

Our results and the assumptions driving them naturally suggest several directions for future research.

First, we assumed that externalities are homogeneous within a profession. However, in reality, externalities are highly heterogeneous within professions. As Mankiw and Whinston (1986) emphasize, entrepreneurial firm formation may be excessively supplied if firms are simply imitating the products of existing firms just as Hirshleifer (1971) emphasized that high-speed trading is over-supplied, while Posner and Weyl (2013) show that long-term price discovery of large bubbles is just as likely to be undersupplied as are innovative breakthroughs. Thus many of the largest gains may come from reallocations within a profession between productive and unproductive activities. Uniform income taxation is unlikely to be a sufficient tool to achieve this reallocation. Mechanisms that do are an exciting direction for future research.

Second, while the optimality of softening incentives to promote efficient allocation of

labor is, to our knowledge, largely unexplored in optimal tax theory, it is widely understood in corporate finance and agency theory.¹⁷ In particular, Holmström and Milgrom (1991) argue that strong material incentives for one, observable dimension of work effort may reduce effort along unobservable dimensions to the extent that effort along the two dimensions is substitutable. Our theory is similar, except that the substitutability arises from the (un-modeled and empirically calibrated) correlation between income levels and the externalities of professions. This alternative micro-foundation of their model also acts as a basis for calibrating it empirically in a taxation context; we are not aware of any analogous calibration of the optimal contract schedule in the original agency context in which they proposed their model.¹⁸ If the productivity, the private payoff from and the allocation of time to tasks could be identified within firms, the theory could be applied with data on the allocation of time to these tasks.¹⁹

Finally, we assumed that profession-specific taxation was, for a variety of reasons, infeasible. However we found that, under the calibrations from the literature, the overwhelming majority of gains possible from first-best, profession-specific taxation could not be achieved by non-discriminatory taxation. This suggests significant welfare could be gained through more targeted instruments and that these merit further investigation. Some of these may even be politically plausible, such as differential subsidies to different types of education, sector-specific output taxation or grants to support work in certain sectors such as research and development.

¹⁷A notable exception is the work of Piketty *et al.* (2014), who argue that taxation may reduce effort expended bargaining for higher compensation. While we agree with the spirit of this result and are motivated by the macroeconomic evidence the authors present, we believe that compensation bargaining has a small elasticity compared to career choice and a career-based theory offers a more useful basis for calibration.

¹⁸Slade (1996) tests directional comparative statics of the Holmström and Milgrom model, but does not structurally calibrate an optimal contract. As far as we know the only other paper to exploit the equivalence of agency and tax theory to link of quantitative optimal tax work to the largely theoretical agency literature is Prendergast (2013), albeit in the more classical context of the Vickrey insurance-incentives trade-off.

¹⁹For example, a friend of one of the authors who works at an investment bank reports that, “I spend one-third of my time creating profits for the firm, one-third of my time ensuring I get credit for those profits and one-third of my time ensuring that I get paid for the profits I got credit for.” Whatever the actual proportions, presumably optimal compensation structure is highly sensitive to these proportions.

References

- ABREU, D. and BRUNNERMEIER, M. K. (2003). Bubbles and Crashes. *Econometrica*, **71** (1), 173–204.
- ACEMOGLU, D. (1995). Reward structures and the allocation of talent. *European Economic Review*, **39** (1), 17–33.
- ALLEN, F., MORRIS, S. and SHIN, H. S. (2006). Beauty Contests and Iterated Expectations in Asset Markets. *Review of Financial Studies*, **19** (3), 719–752.
- ALONSO, W. (1964). *Location and Land Use*. Cambridge, MA: Harvard University Press.
- ATKINSON, A. B., PIKETTY, T. and SAEZ, E. (2011). Top incomes in the long run of history. *Journal of Economic Literature*, **49** (1), 3–71.
- BAGWELL, K. (2007). The economic analysis of advertising. In M. Armstrong and R. H. Porter (eds.), *Handbook of Industrial Organization*, vol. 3, Amsterdam: North-Holland.
- BAI, J., PHILIPPON, T. and SAVOV, A. (2013). Have financial markets become more informative?, <http://pages.stern.nyu.edu/~tphilipp/papers/BaiPhilipponSavov.pdf>.
- BAKIJ, J., COLE, A. and HEIM, B. T. (2012). Job and income growth of top earners and the causes of changing income inequality: Evidence from u.s. tax return data, <http://web.williams.edu/Economics/wp/BakijaColeHeimJobsIncomeGrowthTopEarners.pdf>.
- BAUMOL, W. J. (1990). Entrepreneurship: Productive, unproductive and destructive. *Journal of Political Economy*, **98** (5), 893–921.
- BAYER, P., GEISLER, C. and ROBERTS, J. W. (2013). Speculators and Middlemen: The Role of Intermediaries in the Housing Market. *Working Paper*.
- BECKER, G. S. and MURPHY, K. M. (1993). A simple theory of advertising as a good or a bad. *Quarterly Journal of Economics*, **108** (4), 941–964.
- BERTRAND, M. and MULLAINATHAN, S. (2001). Are ceos rewarded for luck? the ones without principals are. *Quarterly Journal of Economics*, **116** (3), 901–932.
- BLOOM, N., EIFERT, B., MAHAJAN, A., MCKENZIE, D. and ROBERTS, J. (2013). Does management matter? evidence from india. *Quarterly Journal of Economics*, **128** (1), 1–51.

- , PROPPER, C., SELLER, S. and REENEN, J. V. (2011). The impact of competition on management quality: Evidence from public hospitals, http://www.stanford.edu/~nbloom/Hospitals_2011.pdf.
- , SADUN, R. and REENEN, J. V. (2012). Americans do it better: Us multinationals and the productivity miracle. *American Economic Review*, **102** (1), 167–201.
- BROWN, P. (2011). American voters split on government shutdown. *Quinnipiac University Poll*.
- BURCHFIELD, M., OVERMAN, H. G., PUGA, D. and TURNER, M. A. (2006). Causes of Sprawl: A Portrait from Space. *Quarterly Journal of Economics*, **121** (2), 587–633.
- BURNSIDE, C., EICHENBAUM, M. and REBELO, S. (2013). Understanding Booms and Busts in Housing Markets. *Working Paper*.
- BYRD, R. H., NOCEDAL, J. and WALTZ, R. A. (2006). Knitro: An integrated package for nonlinear optimization. In G. D. Pillo and M. Roma (eds.), *Large-Scale Nonlinear Optimization*, New York: Springer Science+Business Media.
- CASE, K. E. and SHILLER, R. J. (1989). The efficiency of the market for single-family homes. *The American Economic Review*, **79** (1), 125–137.
- , — and THOMPSON, A. K. (2012). What Have They Been Thinking? Homebuyer Behavior in Hot and Cold Markets. *Brookings Papers on Economic Activity*, pp. 265–298.
- CHEN, J., HONG, H. and STEIN, J. C. (2002). Breadth of Ownership and Stock Returns. *Journal of Financial Economics*, **66**, 171–205.
- CHENG, I.-H., RAINA, S. and XIONG, W. (2014). Wall street and the housing bubble. *American Economic Review*, **Forthcoming**.
- CHETTY, R., FRIEDMAN, J. and ROCKOFF, J. E. (2013a). Measuring the impact of teachers i: Evaluating bias in teacher value-added estimates, <http://obs.rc.fas.harvard.edu/chetty/w19423.pdf>.
- , — and — (2013b). Measuring the impact of teachers ii: Teacher value-added and student outcomes, <http://obs.rc.fas.harvard.edu/chetty/w19424.pdf>.
- CLARK, M. A., CHIANG, H. S., SILVA, T., MCCONNELL, S., SONNENFELD, K., ERBE, A. and PUMA, M. (2013). *The Effectiveness of Secondary Math Teachers from Teach for America and the Teaching Fellows Program*. Tech. rep., National Center for Education Evaluation and Regional Assistance, Institute of Education Sciences, U. S. Department of Education.
- CNN/ORC (2011). Wall street attitudes poll, <http://i2.cdn.turner.com/cnn/2011/images/10/24/rel17e.pdf>.
- COLOMBI, R. (1990). A new model of income distributions: The pareto lognormal distribution. In C. Dagum and M. Zenga (eds.), *Income and Wealth Distribution, Inequality and Poverty*, Berlin: Springer, pp. 18–32.

- CUTLER, D. M., POTERBA, J. M. and SUMMERS, L. H. (1991). Speculative dynamics. *The Review of Economic Studies*, **58** (3), 529–546.
- DAVIS, M. A. and HEATHCOTE, J. (2007). The Price and Quantity of Residential Land in the United States. *Journal of Monetary Economics*, **54** (8), 2595–2620.
- D’AVOLIO, G. (2002). The market for borrowing stock. *Journal of Financial Economics*, **66** (2), 271–306.
- DEFUSCO, A., DING, W., FERREIRA, F. and GYOURKO, J. (2013). The role of contagion in the last american housing cycle. *Working Paper*.
- DIAMOND, P. and SAEZ, E. (2011). The case for a progressive tax: From basic research to policy recommendations. *Journal of Economic Perspectives*, **25** (4), 165–190.
- DIAMOND, P. A. (1998). Optimal income taxation: An example with a u-shaped pattern of optimal marginal tax rates. *American Economic Review*, **88** (1), 83–95.
- DIXIT, A. and NORMAN, V. (1978). Advertising and welfare. *Bell Journal of Economics*, **9** (1), 1–17.
- EDMANS, A. and GABAIX, X. (2009). Is ceo pay really inefficient? a survey of new optimal contracting theories. *European Financial Management*, **15** (3), 486–496.
- FAMA, E. F. (1970). Efficient Capital Markets: A Review of Theory and Empirical Work. *Journal of Finance*, **25** (2), 383–417.
- FARHI, E. and WERNING, I. (2012). Capital taxation: Quantitative explorations of the inverse euler equation. *Journal of Political Economy*, **120** (3), 398–445.
- FLORIDA, R. and MELLANDER, C. (2010). There goes the metro: How and why bohemians, artists and gays affect regional housing values. *Journal of Economic Geography*, **10** (2), 167–188.
- FRENCH, K. R. (2008). The cost of active investing. *Journal of Finance*, **63** (4), 1537–1573.
- GABAIX, X. and LANDIER, A. (2008). Why has ceo pay increased so much? *Quarterly Journal of Economics*, **123** (1), 49–100.
- GAO, Z. (2014). Housing boom and bust with elastic supplies. *Working Paper*.
- GEANAKOPLOS, J. (2009). The Leverage Cycle. *NBER Macroeconomic Annual*, pp. 1–65.
- GLAESER, E. L. (2013). A Nation of Gamblers: Real Estate Speculation and American History. *American Economic Review*, **103** (3), 1–42.
- and GYOURKO, J. (2005). Urban decline and durable housing. *Journal of Political Economy*, **113** (2).
- and — (2007). Arbitrage in housing markets. *NBER Working Papers*.

- , — and SAIZ, A. (2008). Housing supply and housing bubbles. *Journal of Urban Economics*, **64** (2), 198–217.
- , — and SAKS, R. E. (2005). Why have housing prices gone up? *American Economic Review*, pp. 329–333.
- and KAHN, M. E. (2004). Sprawl and Urban Growth. *Handbook of Urban Economics*, **4**.
- and KOHLHASE, J. E. (2004). Cities, regions, and the decline of transport costs. *Papers in Regional Science*, **83**, 197–228.
- GOLDIN, C., KATZ, L. F., HAUSMAN, N. and WARD, B. (2013). *Harvard and Beyond Project*. Survey data project, Harvard University.
- GRENADIER, S. R. (1996). The strategic exercise of options: Development cascades and overbuilding in real estate markets. *The Journal of Finance*, **51** (5), 1653–1679.
- GROSSMAN, S. (1976). On the efficiency of competitive stock markets where trades have diverse information. *The Journal of Finance*, **31** (2), 573–585.
- GUERRIERI, V., HARTLEY, D. and HURST, E. (2013). Endogenous Gentrification and Housing Price Dynamics. *Working Paper*.
- GURUN, U. G., MATVOX, G. and SERU, A. (2013). Advertising expensive mortgages, http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2233380.
- GYOURKO, J. (2009). Housing Supply. *Annual Review of Economics*, **1**, 295–318.
- , MAYER, C. and SINAI, T. (2013). Superstar cities. *American Economic Journal: Economic Policy*, **5** (4), 167–199.
- and SAIZ, A. (2006). Construction Costs and the Supply of Housing Structure. *Journal of Regional Science*, **46** (4), 661–680.
- HACKER, J. S. and PIERSON, P. (2010). *Winner-Take-All Politics*. New York: Simon & Schuster.
- HANSEN, L. P. (1982). Large sample properties of generalized method of moments estimators. *Econometrica*, pp. 1029–1054.
- HARRISON, J. M. and KREPS, D. M. (1978). Speculative investor behavior in a stock market with heterogeneous expectations. *The Quarterly Journal of Economics*, **92** (2), 323.
- HAUSMAN, J. (2012). Contingent valuation: From dubious to hopeless. *Journal of Economic Perspectives*, **26** (4), 43–56.
- HELLMANZIK, C. (2010). Location matters: Estimating cluster premiums for prominent modern artists. *European Economic Review*, **54** (2), 199–218.
- HELLWIG, M. F. (1980). On the aggregation of information in competitive markets. *Journal of Economic Theory*, **22** (3).

- HENDERSON, J. V. and IOANNIDES, Y. M. (1983). A Model of Housing Tenure Choice. *American Economic Review*, **73** (1), 98–113.
- HICKS, J. R. (1939). The foundations of welfare economics. *Economic Journal*, **49** (196), 696–712.
- HIRSHLEIFER, J. (1971). The private and social value of information and the reward to inventive activity. *American Economic Review*, **61** (4), 561–574.
- HOBERG, G. and PHILLIPS, G. (2010). Real and Financial Industry Booms and Busts. *Journal of Finance*, **65** (1), 45–86.
- HOLMSTRÖM, B. and MILGROM, P. (1991). Multitask principal-agent analyses: Incentive contracts, asset ownership, and job design. *Journal of Law, Economics & Organization*, **7** (Special Issue), 25–52.
- HONG, H., SCHEINKMAN, J. and XIONG, W. (2006). Asset float and speculative bubbles. *The journal of finance*, **61** (3), 1073–1117.
- and SRAER, D. (2012). Speculative Betas. *Working Paper*.
- JACOBS, J. (1961). *The Death and Life of Great American Cities*. New York: Random House.
- JAFFE, A. B. (1989). Real effects of academic research. *American Economic Review*, **79** (5), 957–970.
- KALDOR, N. (1939). Welfare propositions in economics and interpersonal comparisons of utility. *Economic Journal*, **49** (145), 549–552.
- KAPLAN, S. N. and RAUH, J. D. (2010). Wall street and main street: What contributes to the rise in the highest incomes. *Review of Financial Studies*, **23** (3), 1004–1050.
- KEANE, M. and ROGERSON, R. (2012). Micro and macro labor supply elasticities: A reassessment of conventional wisdom. *Journal of Economic Literature*, **50** (2), 464–476.
- KEANE, M. P. (2011). Labor supply and taxes: A survey. *Journal of Economic Literature*, **49** (4), 961–1075.
- KIRKEBØEN, L., LEUVEN, E. and MOGSTAD, M. (2014). Field of study, earnings, and self-selection, http://isites.harvard.edu/fs/docs/icb.topic1296633.files/fieldofstudy_harvard_mm.pdf.
- Land Advisors (2010). *Metro Phoenix Market Overview: First Quarter 2010*. Land Advisors Organization.
- LANDVOIGT, T. (2011). Housing Demand during the Boom: The Role of Expectations and Credit Constraints. *Working Paper*.
- , PIAZZESI, M. and SCHNEIDER, M. (2013). The Housing Market(s) of San Diego. *Working Paper*.
- LAVY, V. and ABRAMITZKY, R. (Forthcoming). How responsive is investment in schooling to changes in redistributive policies and in returns? *Econometrica*.

- LAZEAR, E. P. (1999). Culture and language. *Journal of Political Economy*, **107** (S6), S95–S126.
- LEAMER, E. (2007). Housing is the business cycle. *NBER Working Paper*, **13428**.
- MALMENDIER, U. and TATE, G. (2009). Superstar ceos. *Quarterly Journal of Economics*, **124** (4), 1593–1638.
- MANKIW, N. G. (2010). Spreading the wealth around: Reflections inspired by joe the plumber. *Eastern Economic Journal*, **36**, 285–298.
- and WHINSTON, M. D. (1986). Free entry and social inefficiency. *The RAND Journal of Economics*, **17** (1), 48–58.
- MARX, K. (1867). *Das Kapital, Kritik der Politischen Ökonomie*. Hamburg, Germany: Verlag von Otto Meissner.
- MAYER, C. J. and SOMERVILLE, C. T. (2000). Land Use Regulation and New Construction. *Regional Science and Urban Economics*, **30** (6), 639–662.
- MCCONNELL, J., LINDROOTH, R., WHOLEY, D. and MADDOX, T. (2013). Management practices and the quality of care in cardiac units. *Journal of the American Medical Association: Internal Medicine*, **173** (8), 684–692.
- MCCORMACK, J., PROPPER, C. and SMITH, S. (Forthcoming). Herding cats? management and university performance. *Economic Journal*.
- McKINLEY, J. and PALMER, G. (2007). Nevada Learns to Cash In on Sales of Federal Land. *New York Times*.
- MIAN, A., RAO, K. and SUFI, A. (2013). Household Balance Sheets, Consumption, and the Economic Slump. *Working Paper*.
- and SUFI, A. (2009). The consequences of mortgage credit expansion: Evidence from the US mortgage default crisis. *The Quarterly Journal of Economics*, **124** (4), 1449.
- and — (2011). House Prices, Home Equity-Based Borrowing, and the US Household Leverage Crisis. *American Economic Review*, **101**, 2132–56.
- MILGROM, P. R. and WEBER, R. J. (1982). A theory of auctions and competitive bidding. *Econometrica*, **50** (5), 1089–1122.
- MILLER, E. M. (1977). Risk, uncertainty, and divergence of opinion. *The Journal of Finance*, **32** (4), 1151–1168.
- MIRRLEES, J. A. (1971). An exploration in the theory of optimal taxation. *Review of Economic Studies*, **38** (2), 175–208.
- MORRIS, S. (1996). Speculative investor behavior and learning. *The Quarterly Journal of Economics*, **111** (4), 1111.
- MURPHY, K. M., SHLEIFER, A. and VISHNY, R. W. (1991). The allocation of talent: Implications for growth. *Quarterly Journal of Economics*, **106** (2), 503–530.

- and TOPEL, R. H. (2006). The value of health and longevity. *Journal of Political Economy*, **114** (5), 871–904.
- NEWKEY, W. K. (1984). A method of moments interpretation of sequential estimators. *Economics Letters*, **14** (2), 201–206.
- and MCFADDEN, D. (1994). Large sample estimation and hypothesis testing. *Handbook of econometrics*, **4**, 2111–2245.
- NICHOLS, J. B., OLINER, S. D. and MULHALL, M. R. (2010). Commercial and residential land prices across the United States. *Finance and Economics Discussion Series*.
- ONSTED, J. A. (2009). *The Effectiveness of the Williamson Act: A Spatial Analysis*. VDM Verlag Dr. Müller.
- PARKER, K. (2012). Yes, the rich are different. *Pew Social and Demographic Trends*, **August 27**.
- PÁSTOR, L. and VERONESI, P. (2003). Stock valuation and learning about profitability. *The Journal of Finance*, **58** (5), 1749–1790.
- and VERONESI, P. (2009). Technological revolutions and stock prices. *The American economic review*, **99** (4), 1451–1483.
- PHILIPPON, T. (2010). Financiers versus engineers: Should the financial sector be taxed or subsidized? *American Economic Journal: Microeconomics*, **2** (3), 158–182.
- (2013). Has the u.s. finance industry become less efficient? on the theory and measurement of financial intermediation, http://pages.stern.nyu.edu/~tphilipp/papers/Finance_Efficiency.pdf.
- and RESHEF, A. (2012). Wages and human capital in the u.s. finance industry: 1909–2006. *Quarterly Journal of Economics*, **127** (4), 1551–1609.
- PIAZZESI, M. and SCHNEIDER, M. (2009). Momentum traders in the housing market: Survey evidence and a search model. *The American economic review*, **99** (2), 406–411.
- PIKETTY, T. and SAEZ, E. (2003). Income inequality in the united states, 1913–1998. *Quarterly Journal of Economics*, **118** (1), 1–41.
- , — and STANTCHEVA, S. (2014). Optimal taxation of top labor incomes: A tale of three elasticities. *American Economic Journal: Economic Policy*, **6** (1), 230–271.
- POLINSKY, A. M. and ELLWOOD, D. T. (1979). An empirical reconciliation of micro and grouped estimates of the demand for housing. *The Review of Economics and Statistics*, **61** (2), 199–205.
- POSNER, E. and WEYL, E. G. (2013). Benefit-cost analysis for financial regulation. *American Economic Review Papers and Proceedings*, **103** (3).
- POTERBA, J. M. (1984). Tax subsidies to owner-occupied housing: an asset-market approach. *The Quarterly Journal of Economics*, **99** (4), 729.

- PRENDERGAST, C. (2013). The empirical content of pay-for-performance, <http://businessinnovation.berkeley.edu/WilliamsonSeminar/canice-empirical%20content022813.pdf>.
- RAND, A. (1957). *Atlas Shrugged*. New York: Random House.
- RIZZO, J. A. (1999). Advertising and competition in the ethical pharmaceutical industry: The case of antihypertensive drug. *Journal of Law and Economics*, **42** (1), 89–116.
- ROBACK, J. (1982). Wages, rents, and the quality of life. *The Journal of Political Economy*, **90** (6), 1257–1278.
- ROSEN, S. (1979). Wage-based indexes of urban quality of life. *Current issues in urban economics*, **3**.
- (1981). The economics of superstars. *American Economic Review*, **71** (5), 845–858.
- ROTHSCHILD, C. and SCHEUER, F. (2014a). Optimal taxation with rent-seeking, http://www.stanford.edu/~scheuer/rent_seeking.pdf.
- and — (2014b). A theory of income taxation under multidimensional skill heterogeneity, <http://www.stanford.edu/~scheuer/multidim.pdf>.
- RTCSNV (2012). *Regional Transportation Plan 2013-2035*. Regional Transportation Commission of Southern Nevada.
- SAEZ, E. (2001). Using elasticities to derive optimal income tax rates. *The Review of Economic Studies*, **68** (1), 205–229.
- , SLEMROD, J. and GIERTZ, S. H. (2012). The elasticity of taxable income with respect to marginal tax rates: A critical review. *Journal of Economic Literature*, **50** (1), 3–50.
- SAIZ, A. (2003). Room in the kitchen for the melting pot: Immigration and rental prices. *Review of Economics and Statistics*, **85** (3), 502–521.
- (2010). The Geographic Determinants of Housing Supply. *Quarterly Journal of Economics*, **125** (3), 1253–1296.
- SCHEINKMAN, J. A. and XIONG, W. (2003). Overconfidence and speculative bubbles. *Journal of Political Economy*, **111** (6), 1183–1220.
- SEGEL, A. I., FELDMAN, G. S., LIU, J. T. and WILLIAMSON, E. C. (2011). *Stuyvesant Town - Peter Cooper Village: America's Largest Foreclosure*. Harvard Business School.
- SHILLER, R. J. (2005). *Irrational Exuberance*. Princeton University Press.
- SIMSEK, A. (2013a). Belief Disagreements and Collateral Constraints. *Econometrica*, **81** (1), 1–53.
- (2013b). Speculation and Risk-Sharing with New Financial Assets. *Quarterly Journal of Economics*, **128** (3), 1365–1396.

- SKOCPOL, T. and WILLIAMSON, V. (2012). *The Tea Party and the Remaking of Republican Conservatism*. Oxford: Oxford University Press.
- SLADE, M. E. (1996). Multitask agency and contract choice: An empirical exploration. *International Economic Review*, **37** (2), 465–486.
- SMALL, K. A. and ROSEN, H. S. (1981). Applied welfare economics with discrete choice models. *Econometrica*, **49** (1), 105–130.
- SOMERVILLE, C. T. (1999). The Industrial Organization of Housing Supply: Market Activity, Land Supply and the Size of Homebuilder Firms. *Real Estate Economics*, **27** (4), 669–694.
- SOO, C. K. (2013). Quantifying Animal Spirits: News Media and Sentiment in the Housing Market. *Working Paper*.
- SUHER, M. (2013). Future House Price Expectations in the Recent Boom and Bust. *Working Paper*.
- THE TAX FOUNDATION (2013). *U.S. Federal Individual Income Tax Rates History, 1862-2013 (Nominal and Inflation-Adjusted Brackets)*. Tech. rep., Tax Foundation.
- TITMAN, S. (1983). Urban land prices under uncertainty. *American Economic Review*, **75** (3), 505–514.
- TOPEL, R. and ROSEN, S. (1988). Housing Investment in the United States. *Journal of Political Economy*, **96** (4), 718–740.
- VAN NIEUWERBURGH, S. and WEILL, P.-O. (2010). Why has house price dispersion gone up? *The Review of Economic Studies*, **77** (4), 1567–1606.
- VEIGA, A. and WEYL, E. G. (2013). The leibniz rule for multidimensional heterogeneity, http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2344812.
- VICKREY, W. (1945). Measuring marginal utility by reactions to risk. *Econometrica*, **13** (4), 319–333.
- WILLIAMS, H. C. W. L. (9). On the formation of travel demand models and economic evaluation measures of user benefit. *Environment and Planning A*, **3** (285–344).
- WILSON, R. B. (1993). *Nonlinear Pricing*. Oxford: Oxford University Press.

Appendix A

Appendix to Chapter 1

A.1 Estimation Details

A.1.1 Sequential Two-Step GMM Estimator

Let Z_t^i denote a vector of observed variables that correspond to observation i at period t . This vector may include lagged variables. Denote by ζ the vector of structural parameters that we want to estimate. In our model, the parameter vector ζ corresponds to $(\delta, \theta, \sigma, c_1, c_2)$. Specifically, we use ζ to summarize the vector of wage-related housing demand parameters, (δ, θ, σ) , and γ to denote the vector of housing supply parameters (c_1, c_2) .

We split the vector of moment functions provided by the model into a subvector that depends only on the wage-related structural parameters ζ , $f(Z_t^i; \zeta)$, and the remaining subvector of moment functions that depends both on ζ and γ , $v(Z_t^i; \zeta, \gamma)$. Therefore, using the set of moment functions $f(\cdot)$, we can obtain GMM estimates of ζ that do not depend on the value of γ , $\hat{\zeta}^{SEQ}$. Using the vector of moment functions $v(\cdot)$ and our estimates of ζ , we then estimate γ in our second step, $\hat{\gamma}^{SEQ}$. These estimates of γ will depend on the values estimated for ζ in the first step.

We estimate $\hat{\zeta}^{SEQ}$ by minimizing the objective function:

$$\hat{Q}_2(\zeta) = \left[(NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T f(Z_t^i; \zeta) \right]' \hat{W}_{ff} \left[(NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T f(Z_t^i; \zeta) \right].$$

The weighting matrix \widehat{W}_{ff} is defined as

$$\widehat{W}_{ff} = \left[(NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T f(Z_t^i; \hat{\xi}_1) \cdot f(Z_t^i; \hat{\xi}_1)' \right]^{-1},$$

and $\hat{\xi}_1$ minimizes the first stage objective function

$$\widehat{Q}_1(\xi) = \left[(NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T f(Z_t^i; \xi) \right]' I \left[(NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T f(Z_t^i; \xi) \right],$$

where I denotes the identity matrix. Given that this estimate $\hat{\xi}^{SEQ}$ does not depend on the value of γ , we compute its asymptotic variance as

$$Var(\hat{\xi}^{SEQ}) = \left(\widehat{F}_\xi' \widehat{W}_{ff}^{-1} \widehat{F}_\xi \right)^{-1},$$

where \widehat{W}_{ff} is defined above and \widehat{F}_ξ is

$$\widehat{F}_\xi = (NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T \frac{\partial}{\partial \xi} f(Z_t^i; \xi).$$

Using this initial estimate of ξ , we compute an estimate of γ by minimizing the following objective function:

$$\widehat{Q}_2(\gamma; \hat{\xi}^{SEQ}) = \left[(NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T v(Z_t^i; \gamma, \hat{\xi}^{SEQ}) \right]' \widehat{W}_{vv}(\hat{\xi}^{SEQ}) \left[(NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T v(Z_t^i; \gamma, \hat{\xi}^{SEQ}) \right],$$

where $\widehat{W}_{vv}(\hat{\xi}^{SEQ})$ is

$$\widehat{W}_{vv}(\hat{\xi}^{SEQ}) = \left[(NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T v(Z_t^i; \hat{\gamma}_1, \hat{\xi}^{SEQ}) \cdot v(Z_t^i; \hat{\gamma}_1, \hat{\xi}^{SEQ})' \right]^{-1}$$

and $\hat{\gamma}_1$ minimizes the first stage objective function

$$\widehat{Q}_1(\gamma; \hat{\xi}^{SEQ}) = \left[(NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T v(Z_t^i; \gamma, \hat{\xi}^{SEQ}) \right]' I \left[(NT)^{-1} \sum_{i=1}^N \sum_{t=1}^T v(Z_t^i; \gamma, \hat{\xi}^{SEQ}) \right].$$

The correct formula for the asymptotic variance of $\hat{\gamma}^{SEQ}$ must account for the fact that its distribution depends not only on the random vector $\{Z_t^i; \forall i, t\}$ but also on the additional random vector $\hat{\xi}^{SEQ}$. Newey (1984) provides the correct formula for the asymptotic variance

of the second step estimator:

$$\text{Var}(\hat{\gamma}^{SEQ}) = \left[\hat{V}_\gamma' \hat{W}_{vv}^{-1} \hat{V}_\gamma \right]^{-1} + \hat{V}_\gamma^{-1} \hat{V}_\xi \left[\hat{F}_\xi \hat{W}_{ff}^{-1} \hat{F}_\xi' \right]^{-1} \hat{V}_\xi' \hat{V}_\gamma^{-1'} - \hat{V}_\gamma^{-1} \left[\hat{V}_\xi \hat{F}_\xi^{-1} \hat{W}_{fv} + \hat{W}_{vf} \hat{F}_\xi^{-1'} \hat{V}_\xi' \right] \hat{V}_\gamma^{-1'}.$$

Following Newey and McFadden (1994), the sequential GMM estimators belong to the more general family of *extremum* estimators. This guarantees that they are consistent, asymptotically normal, and have the asymptotic variance described above.

A.1.2 Moment Conditions

Estimation of Housing Demand Parameters

The vectorial moment condition

$$\mathbb{E}[f(\tilde{W}^i; (\delta, \theta, \sigma))] = 0$$

is based on the following vector of moment functions:

$$f(\tilde{W}^i; (\delta, \theta, \sigma)) = \begin{cases} \tau_t^i \\ \tau_t^i \tilde{W}_{t-s}^i & \forall s \geq 3 \\ (\tau_t^i)^2 - (2\theta^2 - 2\theta + 2)\sigma_\epsilon^2 \\ \tau_t^i \tau_{t-1}^i - (-\theta^2 + 2\theta - 1)\sigma_\epsilon^2 \\ \tau_t^i \tau_{t-2}^i - (-\theta)\sigma_\epsilon^2 \end{cases},$$

with

$$\tau_t^i = \Delta \tilde{W}_t^i - \delta \Delta \tilde{W}_{t-1}^i - (1 - \delta)w_1^a = \epsilon_t^i + (\theta - 1)\epsilon_{t-1}^i - \theta\epsilon_{t-2}^i,$$

and $\Delta \tilde{W}_t^i = \tilde{W}_t^i - \tilde{W}_{t-1}^i$. Intuitively, one can think of the random variable τ_t^i as close to (but not exactly) a double-difference of the productivity measure \tilde{W} . The moment function $f(\tilde{W}^i; (\delta, \theta, \sigma))$ is based on the expectation, variance, and serial correlation of this double difference, as well as its covariance with lagged values of the productivity measure \tilde{W} .

Estimation of Housing Supply Parameters

The vectorial moment condition

$$\mathbb{E}[v(\tilde{W}^i; (\delta, \theta, \sigma))] = 0$$

is based on the following vector of moment functions:

$$v(H^i, N^i, I^i; (c_1, c_2)) = \begin{cases} \nu_t^i & \\ \nu_t^i N_{t-s}^i & \forall s \geq 1 \\ \kappa_t^i & \\ \kappa_t^i N_{t-s}^i & \forall s \geq 0 \end{cases} ',$$

$$\begin{aligned} & (\nu_t^i)^2 - \left[\frac{(\bar{\phi} + \hat{\theta})^2}{(\bar{\phi} - \hat{\delta})^2} + \hat{\theta}^2 \right] \hat{\sigma}_\epsilon^2 \\ & (\kappa_t^i)^2 - \left[\frac{(1+r)^2 (\hat{\delta} + \hat{\theta})^2}{(c_1)^2 (\bar{\phi} - \hat{\delta})^2} \right] \sigma_\epsilon^2 \end{aligned}$$

with

$$\begin{aligned} \nu_t^i &= ((H_t^i - \hat{H}_t^i) - \hat{\delta}(H_{t-1}^i - \hat{H}_{t-1}^i) + \frac{\alpha(1+r)}{1+r-\phi} ((N_t^i - \hat{N}_t^i) - \hat{\delta}(N_{t-1}^i - \hat{N}_{t-1}^i))), \\ \kappa_t^i &= ((I_t^i - \hat{I}_t^i) - \hat{\delta}(I_{t-1}^i - \hat{I}_{t-1}^i) + (1-\phi)((N_t^i - \hat{N}_t^i) - \hat{\delta}(N_{t-1}^i - \hat{N}_{t-1}^i))). \end{aligned}$$

Intuitively, one can think of the random variables ν and κ as functions of the differences between the current values of the observable variables (H, I, N) and their steady state values, $(\hat{H}, \hat{I}, \hat{N})$. The moment function $v(H^i, N^i, I^i; (c_1, c_2))$ is based on the expectation and variance of ν and κ , as well as their covariances with lagged values of the number of households, N .

A.1.3 Stochastic Processes Predicted by the Model

If shocks are known as they occur, then our model implies the following ARMA(2,3) process for housing prices

$$\Delta H_t^i = a_0^i + a_1 \Delta H_{t-1}^i + a_2 \Delta H_{t-2}^i + b_0 \epsilon_t^i + b_1 \epsilon_{t-1}^i + b_2 \epsilon_{t-2}^i + b_3 \epsilon_{t-3}^i,$$

where a_0^i denotes a metropolitan area effect, and the parameter vector $(a_1, a_2, b_0, b_1, b_2, b_3)$ is restricted in the following way:

$$\begin{aligned}
a_1 &= \phi + \delta, \\
a_2 &= -\phi\delta, \\
b_0 &= \frac{\bar{\phi} + \theta}{\bar{\phi} - \delta}, \\
b_1 &= \frac{\delta + r(\delta + \theta) - \theta(\delta + \phi) - \bar{\phi}(1 + \delta + \phi)}{\bar{\phi} - \delta}, \\
b_2 &= \frac{\phi\bar{\phi} - \theta(1 + r + \phi(\bar{\phi} - 1)) + \delta(\bar{\phi} - 1 - r + \theta + \theta\phi)}{\bar{\phi} - \delta}, \\
b_3 &= \phi\theta.
\end{aligned}$$

The model also predicts an ARMA(2,1) process for the construction of new houses:

$$I_t^i = d_0^i + d_1 I_{t-1}^i + d_2 I_{t-2}^i + e_0 \epsilon_{t-1}^i + e_1 \epsilon_{t-2}^i,$$

where d_0^i denotes a metropolitan area effect and the parameter vector (d_1, d_2, e_1, e_2) is restricted in the following way:

$$\begin{aligned}
d_1 &= \phi + \delta, \\
d_2 &= -\phi\delta, \\
e_0 &= \frac{(1+r)(\delta + \theta)}{c_1(\bar{\phi} - \delta)}, \\
e_1 &= -\frac{(1+r)(\delta + \theta)}{c_1(\bar{\phi} - \delta)}.
\end{aligned}$$

A.2 Definitions of Trend Variables

We write the housing demand equation as

$$H_t - \frac{E(H_{t+1})}{1+r} - \frac{rC}{1+r} = \bar{x} + qt + x_t - \alpha N_t$$

and the housing supply equation as

$$E(H_{t+1}) = C + c_0 t + c_1 I_t + c_2 N_t.$$

Our third equation is the relationship between city population and construction:

$$N_{t+1} = N_t + I_t.$$

We define \hat{H}_t , \hat{N}_t , and \hat{I}_t to be the unique solutions to the non-stochastic elements of these three equations, which are linear. They are the unique solutions to

$$\begin{aligned}\hat{H}_t - \frac{\hat{H}_{t+1}}{1+r} - \frac{rC}{1+r} &= \bar{x} + qt - \alpha\hat{N}_t \\ \hat{H}_{t+1} &= C + c_0t + c_1\hat{I}_t + c_2\hat{N}_t \\ \hat{N}_{t+1} &= \hat{N}_t + \hat{I}_t\end{aligned}$$

and are therefore given by

$$\begin{aligned}\hat{H}_t &= \frac{c_2^2(\bar{x} + r(C + \bar{x})) + \alpha(1+r)(c_2C - c_0c_1)}{c_2(rc_2 + \alpha(1+r))} + \\ &\quad \frac{(1+r)(\alpha c_0 + qc_2)(c_2^2 + \alpha(1+r)(c_1 - c_2))}{c_2(rc_2 + \alpha(1+r))^2} + \\ &\quad \frac{(1+r)(\alpha c_0 + qc_2)}{rc_2 + \alpha(1+r)}t; \\ \hat{N}_t &= \frac{rc_0c_1 + (1+r)c_2\bar{x}}{c_2(rc_2 + \alpha(1+r))} - \frac{(1+r)(r(c_1 - c_2) - c_2)(\alpha c_0 + qc_2)}{c_2(rc_2 + \alpha(1+r))^2} + \\ &\quad \frac{q(1+r) - rc_0}{rc_2 + \alpha(1+r)}t; \\ \hat{I}_t &= \frac{q(1+r) - rc_0}{rc_2 + \alpha(1+r)}.\end{aligned}$$

A.3 Proofs

A.3.1 Proof of Lemma 1

Let $h = H - \hat{H}$, $n = N - \hat{N}$, and $i = I - \hat{I}$, where H, N, I constitute a solution to the key equations. We must show that $\lim_{j \rightarrow \infty} E(h_{t+j}) = E(n_{t+j}) = E(i_{t+j}) = 0$. Note that because

\hat{H}_t is linear in t ,

$$\lim_{j \rightarrow \infty} \frac{\hat{H}_{t+j}}{(1+r)^j} = 0.$$

Therefore, the transversality condition on H guarantees that

$$\lim_{j \rightarrow \infty} \frac{E(h_{t+j})}{(1+r)^j} = 0.$$

The three key equations reduce to

$$\begin{aligned} h_t - \frac{E(h_{t+1})}{1+r} &= x_t - \alpha n_t \\ E(h_{t+1}) &= c_1 i_t + c_2 n_t \\ n_{t+1} &= n_t + i_t. \end{aligned}$$

Note that because $E(x_{t+j}) = \delta^{j-1} E(x_{t+1})$, $\lim_{j \rightarrow \infty} E(x_{t+j}) = 0$. We can therefore deduce from the first key equation that

$$\lim_{j \rightarrow \infty} \frac{E(n_{t+j})}{(1+r)^j} = 0. \quad (\text{A.1})$$

By combining the three key equations, we obtain the following difference equation:

$$\begin{aligned} (1+r)E(x_{t+1}) &= (1+r)(c_2 - c_1)E(n_t) + ((1+r)(\alpha + c_1) + c_1 - c_2)E(n_{t+1}) - c_1 E(n_{t+2}) \\ &= -c_1(\phi L - I)(\bar{\phi} L - I)E(n_{t+2}). \end{aligned}$$

Here L is the lag operator, I is the identity operator, and $\phi < \bar{\phi}$ are the two roots to the characteristic equation

$$0 = -c_1 y^2 + ((1+r)(\alpha + c_1) + c_1 - c_2)y + (1+r)(c_2 - c_1).$$

Because $\lim_{j \rightarrow \infty} E(x_{t+j}) = 0$, one of the following two equations must hold:

$$\begin{aligned} \lim_{j \rightarrow \infty} (\phi L - I)E(n_{t+j}) &= 0 \\ \lim_{j \rightarrow \infty} (\bar{\phi} L - I)E(n_{t+j}) &= 0. \end{aligned}$$

We claim only the first holds. To prove this, we show that $0 < \phi < 1 < 1 + r < \bar{\phi}$. Recall that $\alpha > 0$ and $c_2 < c_1$. Note that

$$\begin{aligned}
\bar{\phi} &= \frac{(1+r)(\alpha + c_1) + c_1 - c_2 + \sqrt{((1+r)(\alpha + c_1) + c_1 - c_2)^2 - 4(1+r)c_1(c_1 - c_2)}}{2c_1} \\
&> \frac{(1+r)(\alpha + c_1) + c_1 - c_2 + \sqrt{((1+r)c_1 + c_1 - c_2)^2 - 4(1+r)c_1(c_1 - c_2)}}{2c_1} \\
&= \frac{(1+r)(\alpha + c_1) + c_1 - c_2 + rc_1 + c_2}{2c_1} \\
&= 1 + r + \frac{\alpha(1+r)}{2c_1} \\
&> 1 + r.
\end{aligned}$$

We also have

$$\phi + \bar{\phi} = \frac{(1+r)(\alpha + c_1) + c_1 - c_2}{c_1} > 0$$

and

$$\phi\bar{\phi} = (1+r) \left(1 - \frac{c_2}{c_1}\right) \in (0, 1+r).$$

Therefore $0 < \phi < 1$. Because of equation (A.1) and the fact that $\bar{\phi} > 1 + r$, the limit

$$\lim_{j \rightarrow \infty} (\bar{\phi}L - I)E(n_{t+j}) = 0$$

cannot hold. Therefore,

$$\lim_{j \rightarrow \infty} (\phi L - I)E(n_{t+j}) = 0$$

holds, and because $0 < \phi < 1$, this limit implies the limit $\lim_{j \rightarrow \infty} E(n_{t+j}) = 0$. The equation $n_{t+1} = n_t + i_t$ shows that $\lim_{j \rightarrow \infty} E(i_{t+j}) = 0$, and then the equation $E(h_{t+1}) = c_1 i_t + c_2 n_t$ shows that $\lim_{j \rightarrow \infty} E(h_{t+j}) = 0$, concluding the proof.

A.3.2 Proof of Lemma 2

We can write the difference equation for $E(n_{t+2})$ as

$$(1+r)E(x_{t+1}) = -c_1(\phi L - I)(\bar{\phi}I - L^{-1})LE(n_{t+2}) = -c_1(\phi L - I)(\bar{\phi}I - L^{-1})n_{t+1}$$

where we have used the fact that n_{t+1} is known at time t . Because $\bar{\phi} > 1 + r$, the operator $\bar{\phi}I - L^{-1}$ is invertible in the space of functions obeying the transversality condition. We can therefore simplify the above equation to

$$(\phi L - I)n_{t+1} = -\frac{1+r}{c_1}\bar{\phi}^{-1}(I - \bar{\phi}^{-1}L^{-1})^{-1}E(x_{t+1}) = -\frac{1+r}{c_1(\bar{\phi} - \delta)}E(x_{t+1}).$$

Therefore

$$i_t = n_{t+1} - n_t = -(\phi L - I)n_{t+1} - (1 - \phi)n_t = \frac{1+r}{c_1(\bar{\phi} - \delta)}E(x_{t+1}) - (1 - \phi)n_t,$$

which proves the second equation of the lemma. For the first equation of the lemma, note that

$$\begin{aligned} h_t &= x_t - \alpha n_t + \frac{E(h_{t+1})}{1+r} \\ &= x_t - \alpha n_t + \frac{c_1}{1+r}i_t + \frac{c_2}{1+r}n_t \\ &= x_t + \frac{E(x_{t+1})}{\bar{\phi} - \delta} + \left(\frac{c_2 - (1 - \phi)c_1}{1+r} - \alpha \right) n_t \\ &= x_t + \frac{E(x_{t+1})}{\bar{\phi} - \delta} - \frac{\alpha(1+r)}{1+r - \phi}n_t. \end{aligned}$$

The last equality comes verifying that the equality

$$\frac{c_2 - (1 - \phi)c_1}{1+r} - \alpha = -\frac{\alpha(1+r)}{1+r - \phi}$$

follows from the polynomial defining ϕ .

A.3.3 Proof of Proposition 1

We have

$$\begin{aligned}
E(n_{t+j}) &= \phi^j n_t + \sum_{k=0}^{j-1} \phi^k (E(n_{t+j-k}) - E(n_{t+j-k-1})) \\
&= \phi^j n_t + \frac{1+r}{c_1(\bar{\phi} - \delta)} \sum_{k=0}^{j-1} \phi^k E(x_{t+j-k}) \\
&= \phi^j n_t + \frac{1+r}{c_1(\bar{\phi} - \delta)} E(x_{t+1}) \sum_{k=0}^{j-1} \phi^k \delta^{j-1-k} \\
&= \phi^j n_t + \frac{1+r}{c_1(\bar{\phi} - \delta)} \frac{\phi^j - \delta^j}{\phi - \delta} E(x_{t+1}).
\end{aligned}$$

Note that $\hat{N}_{t+j} - \hat{N}_t = j\hat{I}$. Therefore

$$E(N_{t+j} - N_t) = j\hat{I} + \frac{1+r}{c_1(\bar{\phi} - \delta)} \frac{\phi^j - \delta^j}{\phi - \delta} E(x_{t+1}) - (1 - \phi^j)n_t.$$

Next, we have

$$\begin{aligned}
E(I_{t+j}) &= E(N_{t+j+1}) - E(N_{t+j}) \\
&= \hat{I} + \frac{1+r}{c_1(\bar{\phi} - \delta)} \left(\frac{\delta^j(1 - \delta) - \phi^j(1 - \phi)}{\phi - \delta} \right) E(x_{t+1}) - \phi^j(1 - \phi)n_t.
\end{aligned}$$

Finally,

$$\begin{aligned}
E(H_{t+j} - H_t) &= \hat{H}_{t+j} - \hat{H}_t + E(h_{t+j}) - h_t \\
&= \hat{H}_{t+j} - \hat{H}_t + c_1 E(i_{t+j-1}) + c_2 E(n_{t+j}) - h_t,
\end{aligned}$$

and substituting in the previous two results, as well as the formula from Lemma 2, gives the desired equation.

A.3.4 Proof of Proposition 2

Given the hypotheses, we have $x_t = \epsilon_t > 0$ and $E(x_{t+1}) = \delta x_t + \theta \epsilon_t = (\delta + \theta)\epsilon_t > 0$. (The proof is symmetric when $\epsilon_t < 0$.) We also have $n_t = 0$. Therefore, from Lemma 2, prices and construction are both above trend levels at time t . Because $c_2 = 0$, to show that prices

and construction fall below trend at $t + j$ for large j , we must by Proposition 1 show that

$$\frac{\delta^j(1-\delta) - \phi^j(1-\phi)}{\phi - \delta} = \delta^j \frac{1-\delta - (1-\phi)(\phi/\delta)^j}{\phi - \delta}$$

is negative for large j . If $\phi > \delta$, then the numerator is negative for large j , while the denominator is positive. The opposite occurs when $\phi < \delta$. In either case, the fraction is negative for large j , which is what we wanted to show.

A.3.5 Proof of Proposition 3

From Lemma 2, we have

$$I_0 = \hat{I} + \frac{1+r}{c_1(\bar{\phi} - \delta)} \delta \epsilon_0.$$

Therefore

$$N_1 = \hat{N}_1 + \frac{1+r}{c_1(\bar{\phi} - \delta)} \delta \epsilon_0.$$

It follows again from Lemma 2 that

$$I_1 = \hat{I} + \frac{1+r}{c_1(\bar{\phi} - \delta)} \delta (\epsilon_1 + \delta \epsilon_0) - (1-\phi) \frac{1+r}{c_1(\bar{\phi} - \delta)} \delta \epsilon_0.$$

Recall that

$$\hat{I} = \frac{q(1+r) - rc_0}{rc_2 + \alpha(1+r)}.$$

Therefore

$$\text{Cov}(I_0, I_1) = \left(\frac{1+r}{rc_2 + \alpha(1+r)} \right)^2 \text{Var}(q) + \left(\frac{\delta(1+r)}{c_1(\bar{\phi} - \delta)} \right)^2 (\delta + \phi - 1) \text{Var}(\epsilon_0),$$

which is positive as long as

$$\frac{\text{Var}(q)}{\text{Var}(\epsilon_0)} > (1 - \delta - \phi) \left(\frac{\delta(rc_2 + \alpha(1+r))}{c_1(\bar{\phi} - \delta)} \right)^2.$$

Because $\delta > 1 - \phi$ by assumption, this inequality must hold.

We now turn to price growth. From Lemma 2, we have

$$H_0 = \hat{H}_0 + \frac{\bar{\phi}}{\bar{\phi} - \delta} \epsilon_0.$$

We also have

$$H_1 = \hat{H}_0 + \frac{\bar{\phi}}{\bar{\phi} - \delta} \epsilon_1 + \left(\frac{\bar{\phi}}{\bar{\phi} - \delta} - \frac{\alpha(1+r)}{1+r-\phi} \frac{1+r}{c_1(\bar{\phi} - \delta)} \right) \delta \epsilon_0.$$

To obtain H_2 , we first note that

$$N_2 = N_1 + I_1 = \hat{N}_2 + \frac{\delta(1+r)}{c_1(\bar{\phi} - \delta)} (\epsilon_1 + (\delta + \phi - 1)\epsilon_0).$$

Therefore

$$\begin{aligned} H_2 &= \hat{H}_2 + \frac{\bar{\phi}}{\bar{\phi} - \delta} \epsilon_2 + \left(\frac{\bar{\phi}}{\bar{\phi} - \delta} - \frac{\alpha(1+r)}{1+r-\phi} \frac{1+r}{c_1(\bar{\phi} - \delta)} \right) \delta \epsilon_1 \\ &\quad + \left(\frac{\delta \bar{\phi}}{\bar{\phi} - \delta} - \frac{\alpha(1+r)}{1+r-\phi} \frac{(1+r)(\phi + \delta)}{c_1(\bar{\phi} - \delta)} \right) \delta \epsilon_0. \end{aligned}$$

Recall that

$$\hat{H}_{t+1} - \hat{H}_t = \frac{(1+r)(\alpha c_0 + q c_2)}{r c_2 + \alpha(1+r)}.$$

Much algebra produces

$$\text{Cov}(H_2 - H_1, H_1 - H_0) = \left(\frac{(1+r)c_2}{r c_2 + \alpha(1+r)} \right)^2 \text{Var}(q) - \frac{\Omega}{(c_1(\bar{\phi} - \delta))^2} \text{Var}(\epsilon),$$

where $\text{Var}(\epsilon)$ is the common variance of ϵ_0 and ϵ_1 , and Ω is given by

$$\Omega = \left(\frac{\alpha(1+r)^2 \delta}{1+r-\phi} + c_1(1-\delta)\bar{\phi} \right) \left(\frac{(1-\delta-\phi)\alpha(1-r)^2 \delta}{1+r-\phi} + c_1(1-\delta+\delta^2)\bar{\phi} \right).$$

This expression is negative as long as

$$\frac{\text{Var}(q)}{\text{Var}(\epsilon)} < \Omega \left(\frac{r c_2 + \alpha(1+r)}{(1+r)c_1 c_2 (\bar{\phi} - \delta)} \right)^2,$$

which is the bound mentioned in the proposition.

A.4 Calculation of Volatilities in Table 1.1

We showed above that

$$(\phi L - I)n_{t+1} = -\frac{1+r}{c_1(\bar{\phi} - \delta)} E(x_{t+1}).$$

It follows from decomposing n_t into a linear combination of the ϵ terms that

$$\text{Std}(n_t) = \frac{(1+r)(\theta+\delta)}{c_1(\bar{\phi}-\delta)} \sum_{j=1}^{\infty} \frac{\phi^j - \delta^j}{\phi - \delta} \sigma_{\epsilon}.$$

This expression allows us to compute the congestion externality volatility, as it is a constant times this expression.

The other volatility is of a sum of x_t and $E(x_{t+1})$, which can be written as a sum of x_t and ϵ_t . The variance of x is obtained from taking the recursive equation

$$x_t = \delta x_{t-1} + \epsilon_t + \theta \epsilon_{t-1}$$

and then computing the variance of both sides. The result is

$$\sigma_x^2 = \frac{1 + 2\delta\theta + \theta^2}{1 - \delta^2} \sigma_{\epsilon}^2.$$

The covariance of x_t and ϵ_t is just σ_{ϵ}^2 . Therefore, we can readily compute the volatility of the direct wage term, using the standard formula for the variance of a sum of correlated variables.

A.5 BEA Income Data Tables

Table A.1: *Estimated Demand and Supply Parameters: BEA Income Data, 1980-2003*

	Coastal	Sunbelt	Interior
δ	0.80 (0.11)	0.90 (0.08)	0.73 (0.07)
θ	0.16 (0.13)	-0.01 (0.16)	-0.06 (0.13)
σ_ϵ	\$1,200 (200)	\$1,000 (100)	\$800 (80)
Supply			
c_1	6.08 (1.21)	1.00 (0.09)	2.03 (0.35)
c_2	1.88 (0.40)	0.20 (0.03)	0.48 (0.12)

Notes: δ , θ , and σ_ϵ are the autocorrelation parameter, moving average parameter and residual variance of an ARMA(1,1) estimated for the component of wages that is not explained by a linear time trend and a metropolitan area-specific constant. c_1 denotes the derivative of expected future housing prices with respect to current investment in housing construction; and c_2 denote the derivative of the physical capital cost of building a home with respect to the stock of houses. The standard errors for the demand parameters are efficient two-step GMM standard errors. The ones for the supply parameters account for error coming from the demand estimates.

Table A.2: *Volatility and Serial Correlation in House Prices and Construction: BEA Income Data, 1980-2003*

Horizon	Coastal		Sunbelt		Interior	
	Model	Data	Model	Data	Model	Data
<i>Volatility of House Price Changes (\$)</i>						
1 year	5,600	12,650	3,400	2,600	2,300	3,800
3 year	8,800	32,300	5,000	6,500	3,200	9,200
5 year	10,100	44,100	5,600	9,200	3,500	12,600
<i>Serial Correlation of House Price Changes</i>						
1 year	-0.09	0.75	-0.16	0.60	-0.20	0.66
3 year	-0.27	0.09	-0.32	0.21	-0.37	0.17
5 year	-0.36	-0.57	-0.39	-0.24	-0.45	-0.31
<i>Volatility of Construction (units)</i>						
1 year	800	2,600	2,800	5,300	700	2,100
3 year	1,900	6,700	6,700	14,000	1,600	5,100
5 year	2,600	9,800	9,500	19,600	2,200	6,800
<i>Serial Correlation of Construction</i>						
1 year	0.49	0.75	0.56	0.79	0.44	0.73
3 year	0.12	0.27	0.26	0.37	0.05	0.22
5 year	-0.12	-0.27	-0.04	-0.20	-0.29	-0.24

Notes: The moments computed from the data allows the mean of housing price changes and construction to vary across metropolitan areas. The moments generated from the model use the estimates in Table 1.2.

Appendix B

Appendix to Chapter 2

B.1 Micro-foundation of owner-occupancy utility

We present a moral hazard framework in which ownership utility matches the specification of (2.2). Our framework follows the spirit of Henderson and Ioannides (1983)'s treatment of tenure choice, in which maintenance frictions lead some residents to own instead of rent.

Residents derive utility from the particular way their house is “customized”: e.g. the color of the walls, the way the lawn is maintained, et cetera. The set of possible customizations is \mathcal{K} . Resident i 's utility from housing is $v(\sum_{k \in \mathcal{K}} a_{i,k} h_k)$, where $a_{i,k} > 0$ is his preference for k and h_k is the quantity of housing customized that way. Individual customization choices are not contractible, but the right to customize one's house is. If a landlord retains the customization rights, then the tenant cannot customize the house, and the null customization $k = 0$ occurs for which $a_{i,0} = 1$ for all residents i . If the tenant holds these rights, he may choose any $k \in \mathcal{K} \setminus \{0\}$.

Moral hazard arises due to a doomsday customization $k = d$. This customization incurs a cost $\eta(h_d)$ to the owner of the house. All residents prefer this customization to all others: $d = \arg \max_{k \in \mathcal{K}} a_{i,k}$. However, the costs of d outweigh the benefits: for all i and h ,

$$v(a_{i,d}h) < \eta(h).$$

The doomsday customization represents the proclivity of residents to damage a house when they do not bear the costs of doing so.

This inequality prevents landlords from ever selling customization rights to tenants. Suppose the landlord sells the rights. Then the tenant chooses his preferred customization, without taking into account the resultant costs, which the landlord bears. The tenant therefore chooses $k = d$. Knowing this, the landlord demands at least $\eta(h)$ for the customization rights. But the most the tenant is willing to pay is $v(a_{i,d}h) - v(h)$, which is less than $\eta(h)$. Therefore they agree not to trade. The landlord keeps the rights, and $k = 0$. The utility from renting is $v(h)$ because $a_{i,0} = 1$.

An owner-occupant chooses the customization, but also bears the costs if he chooses $k = d$. Let $k(i)$ denote the solution to his optimization problem $\max_{k \in \mathcal{K} \setminus \{0\}} v(a_{i,k}h) - \eta(h)\mathbf{1}_{k=d}$. Due to the costliness of the doomsday customization, the resident never chooses it: $k(i) \neq d$. Indeed, if k' is any other customization, then $v(a_{i,k'}h) > v(a_{i,d}h) - \eta(h)$ due to the above inequality. We define $a_i \equiv a_{i,k(i)}$. The utility from owning is $v(a_i h)$. This form corresponds exactly to (2.2).

B.2 Proofs

B.2.1 Proof of Lemma 3

First we prove that construction occurs in each period. Construction occurs at time 0 because the housing stock starts at 0, and the housing demand equation (2.9) is positive. For a contradiction, let $t_1 > 0$ denote the first period in which construction does not occur. Let $t_2 > t_1$ denote the next time construction occurs (t_2 may be infinite).

We now claim that $r_t^h > r_{t_1-1}^h$ for $t_1 \leq t < t_2$. Along the trend growth path, $x = 0$, so no uncertainty exists and by (2.7), a resident rents if and only if $a_i < 1$. Because F_a has full support on \mathbb{R}^+ , some residents must rent. Landlords hence exist in equilibrium, and their arbitrage equation $p_t^h = r_t^h + \beta p_{t+1}^h$ holds. Aggregate housing demand resulting from the

first-order condition (2.6) is

$$D_t^h(r_t^h) = N_t \left(\int_0^1 (v')^{-1}(r_t^h) dF_a + \int_1^\infty (v')^{-1}(r_t^h/a_i) / a_i dF_a \right). \quad (\text{B.1})$$

By assumption, the housing stock and hence housing demand is the same for $t_1 - 1 \leq t < t_2$. Equation (B.1) decreases in r_t^h because $v'' < 0$. Because $x = 0$ and $g > 0$, N_t increases with t . The left side of (B.1) stays constant for $t_1 - 1 \leq t < t_2$ while N_t increases. Therefore, r_t increases for $t_1 - 1 \leq t < t_2$.

Because construction occurs at $t_1 - 1$, we have $p_{t_1-1}^h = p_{t_1-1}^l + K$, which results from zero homebuilder profits. Zero construction at t_1 can only occur when $p_{t_1}^h \leq p_{t_1}^l + K$, from homebuilder profit maximization. The landlord and landowner arbitrage equations at $t_1 - 1$ deliver $r_{t_1-1}^h \geq r_{t_1-1}^l + (1 - \beta)K$. The quantity of undeveloped land stays constant for $t_1 - 1 \leq t < t_2$ and hence r_t^l does as well, because firm land demand D^l does not change over time. Therefore $r_t^h > r_t^l + (1 - \beta)K$ for $t_1 \leq t < t_2$. Then

$$\begin{aligned} p_{t_1}^h &= \sum_{t_1 \leq t < t_2} \beta^{t-t_1} r_t^h + \beta^{t_2-t_1} p_{t_2}^h \\ &> \sum_{t_1 \leq t < t_2} \beta^{t-t_1} (r_t^l + (1 - \beta)K) + \beta^{t_2-t_1} (p_{t_2}^l + K) \\ &= p_t^l + K, \end{aligned}$$

which contradicts the zero construction inequality $p_{t_1}^h \leq p_t^l + K$. This contradiction proves that construction occurs at all times t .

We now show that rents r_t^h increase over time. Because construction occurs at all t , $p_t^h = p_t^l + K$ for all t . Undeveloped land must always exist because perpetual construction occurs. Therefore, landowners are indifferent between holding land until tomorrow or selling it, so $p_t^l = r_t^l + \beta p_{t+1}^l$. Together with the landlord arbitrage equation, this equation gives $r_t^h = p_t^h - \beta p_{t+1}^h = p_t^l + K - \beta(p_{t+1}^l + K) = r_t^l + (1 - \beta)K$. Equilibrium rents are determined by $S - D^l(r_t^h - (1 - \beta)K) = D_t^h(r_t^h)$, where housing demand comes from (B.1). The left side increases in r_t^h , whereas the right side decreases. N_t increases over time, which shifts up D_t^h . Therefore r_t^h increases as well.

Finally, we can show directly that the supply elasticity decreases over time. The elasticity by definition is

$$\begin{aligned}\epsilon_t^S &= \frac{r_t^h (D^l)' (r_t^h - (1 - \beta)K)}{S - D^l(r_t^h - (1 - \beta)K)} \\ &= \frac{r_t^h}{r_t^h - (1 - \beta)K} \frac{D^l(r_t^h)}{D_t^h(r_t^h)} \frac{r_t^l (D^l)'(r_t^l)}{D^l(r_t^l)} \\ &= \frac{r_t^h}{r_t^h - (1 - \beta)K} \left(\frac{S}{H_t} - 1 \right) \epsilon^l,\end{aligned}$$

which coincides with (2.4). We have shown directly that H_t and r_t^h increase over time. Therefore, when ϵ^l is constant, ϵ_t^S decreases over time.

B.2.2 Proof of Proposition 5

We use (2.5) to write $p_0^h = r_0^h + \beta \tilde{\mathbf{E}} p_1^h$. Let $\partial/\partial x$ denote the partial derivative in which N_0 stays constant. Then $\partial p_0^h/\partial x = \partial r_0^h/\partial x + \partial \beta \tilde{\mathbf{E}} p_1^h/\partial x$. We calculate $\partial r_0^h/\partial x$ by differentiating (2.8) at $x = 0$. Let $d(\cdot) = (v')^{-1}(\cdot)$, and let $b_i = 1 + \beta(\tilde{\mathbf{E}} p_1^h - \mathbf{E}_i p_1^h)/r_0^h$. Note that when $x = 0$, $b_i = 1$ for all i . Then

$$\begin{aligned}-(D^l)' \frac{\partial r_0^h}{\partial x} &= N_0 \int_M \int_0^1 d'(r_0^h) \frac{\partial r_0^h}{\partial x} dF_a dF_\mu + N_0 \int_M d(r_0^h) \frac{\partial b_i}{\partial x} dF_\mu \\ &\quad + N_0 \int_M \int_1^\infty a_i^{-2} d'(r_0^h/a_i) \left(\frac{\partial r_0^h}{\partial x} + \frac{\partial \beta \tilde{\mathbf{E}} p_1^h}{\partial x} - \frac{\partial \beta \mathbf{E}_i p_1^h}{\partial x} \right) dF_a dF_\mu - N_0 \int_M d(r_0^h) \frac{\partial b_i}{\partial x} dF_\mu.\end{aligned}$$

The extensive margins terms for the rental and owner-occupied populations cancel. We simplify this equation to

$$\frac{\partial r_0^h}{\partial x} = - \frac{N_0 \int_M \int_1^\infty a_i^{-2} d'(r_0^h/a_i) \left(\partial \beta \tilde{\mathbf{E}} p_1^h/\partial x - \partial \beta \mathbf{E}_i p_1^h/\partial x \right) dF_a dF_\mu}{(D^l)' + N_0 \int_M \int_0^1 d'(r_0^h) dF_a dF_\mu + N_0 \int_M \int_1^\infty a_i^{-2} d'(r_0^h/a_i) dF_a dF_\mu}.$$

The proposition assumes a constant elasticity of housing demand ϵ^D . This property occurs when individual demand $d(\cdot)$ displays the same constant elasticity. Indeed, from (B.1), the elasticity of housing demand when $x = 0$ is

$$\epsilon^D = - \frac{\int_0^1 r d'(r) dF_a + \int_1^\infty r a_i^{-2} d'(r/a_i) dF_a}{\int_0^1 d(r) dF_a + \int_1^\infty a_i^{-1} d(r/a_i) dF_a},$$

which holds when $rd'(r)/d(r) = -\epsilon^D$ for all r . We can therefore rewrite $\partial r_0^h/\partial x$ as

$$\frac{\partial r_0^h}{\partial x} = -\frac{\epsilon^D N_0 \int_M \int_1^\infty a_i^{-1} d(r_0^h/a_i) \left(\partial \beta \tilde{\mathbf{E}} p_1^h / \partial x - \partial \beta \mathbf{E}_i p_1^h / \partial x \right) dF_a dF_\mu}{r_0^h (D^l)' + \epsilon^D N_0 \int_M \int_0^1 d(r_0^h) dF_a dF_\mu + \epsilon^D N_0 \int_M \int_1^\infty a_i^{-1} d(r_0^h/a_i) dF_a dF_\mu}.$$

Because F_a and F_μ are independent, we can write

$$\begin{aligned} \int_M \int_1^\infty a_i^{-1} d(r_0^h/a_i) \mathbf{E}_i p_1^h dF_a dF_\mu &= \int_1^\infty a_i^{-1} d(r_0^h/a_i) dF_a \int_M \mathbf{E}_i p_1^h dF_\mu \\ &= \int_1^\infty a_i^{-1} d(r_0^h/a_i) dF_a \bar{\mathbf{E}} p_1^h, \end{aligned}$$

where $\bar{\mathbf{E}} p_1^h \equiv \int_M \mathbf{E}_i p_1^h dF_\mu$ is the average belief about p_1^h . Recall from (B.1) that $(h_{i,0}^{rent})^* = d(r_0^h)$ if $a_i < 1$ (and 0 otherwise) and $(h_{i,0}^{own})^* = d(r_0^h/a_i)/a_i$ if $a_i \geq 1$ (and 0 otherwise). The share of housing that is owner-occupied is $\chi = \int_1^\infty a_i^{-1} d(r_0^h/a_i) dF_a / (\int_0^1 d(r_0^h) dF_a + \int_1^\infty a_i^{-1} d(r_0^h/a_i) dF_a)$. We can therefore divide through the equation for $\partial r_0^h/\partial x$ by the total housing stock to get

$$\frac{\partial r_0^h}{\partial x} = -\frac{\epsilon^D \chi \left(\partial \beta \tilde{\mathbf{E}} p_1^h / \partial x - \partial \beta \bar{\mathbf{E}} p_1^h / \partial x \right)}{\epsilon_0^S + \epsilon^D}.$$

Substituting into $\partial p_0^h/\partial x = \partial r_0^h/\partial x + \partial \beta \tilde{\mathbf{E}} p_1^h/\partial x$ yields (2.10) of the proposition.

B.2.3 Proof of Proposition 6

We will calculate the effect of the shock z_t on r_t^h by differentiating the equation $S - D^l(r_t^h - (1 - \beta)K) = D_t^h(r_t^h)$ with respect to x at $x = 0$, where $D_t^h(r_t^h)$ is given by (B.1). This derivative is valid if and only if this equilibrium condition holds for x around 0. The condition holds as long as construction occurs at t . Our first task is thus proving the existence of an open set $I \in \mathbb{R}$ such that $0 \in I$ and for $x \in I$, construction occurs for all t .

As in the proof of Lemma 3, we can prove that construction must occur at t_1 if, conditional on the absence of construction at t_1 , $r_t > r_{t_1-1}^h$ for $t_1 \leq t < t_2$ where t_2 is the next time construction occurs. The key step in that proof was that N_t increases with t . We define an open set I_1 containing 0 such that N_t still increases in t for $x \in I_1$. Because M is uniformly bounded, there exist μ^{min} and μ^{max} such that $\mu^{min} \leq \mu' \leq \mu^{max}$ for all μ' that are coordinates of vectors in M . Recall that $N_{t+1}/N_t = e^{g + (\mu_{t+1} - \mu_t)x}$. Because $g > 0$, the set

$I_1 = (-g/(\mu^{max} - \mu^{min}), g/(\mu^{max} - \mu^{min}))$ is open. For any $x \in I_1$, $N_{t+1}/N_t > 1$. With this result, the proof of this increasing rent condition matches verbatim the proof given in the proof of Lemma 3 when $t_1 > 1$. When $t_1 = 1$, $D_{t_1-1}^h$ is no longer given by (B.1) but instead by (2.9).

The only new fact we must show therefore is that if construction fails to occur at $t = 1$, then $r_0^h < r_1^h$. To do this, we first show that $\tilde{\mathbb{E}}p_1^h - \mathbb{E}_i p_1^h = O(x)$ as $x \rightarrow 0$ for all i . We have $p_1^h = \sum_{t=1}^{\infty} \beta^{t-1} r_t^h$. All residents agree on H_0 and N_0 because they are observable at $t = 0$. Let t_2 be the next time construction occurs given H_0 . Once it occurs it will occur afterward forever due to the arguments in the proof of Lemma 3. In principle residents could disagree about t_2 , but we will now show that for x small enough they do not. While construction does not occur, rents are determined by $H_0 = N_t D_t^h(r_t^h)$ and $S - H_0 = D^l(r_t^l)$. Because N_t increases over time, r_t^h must as well. When construction occurs next period but not today at t , $p_t^h < p_t^l + K$ while $p_{t+1}^h = p_{t+1}^l + K$, so using the landlord and landowner arbitrage equations defined in the proof of Lemma 3, we find that $(D_t^h)^{-1}(H_0/N_t) < (D^l)^{-1}(S - H_0) + (1 - \beta)K$ while construction fails to occur. The first time construction does occur, $t = t_2$, is defined as the lowest value of t for which this inequality fails to hold. Because we are in discrete time, and because the relationships $N_t = N_0 e^{g t + (\mu_t - 1)x}$ and $\mu^{min} \leq \mu_t \leq \mu^{max}$ hold, there exists an open $I_2 \ni 0$ such that when $x \in I_2$, t_2 is the same for all realizations of $\mu \in M$. For $1 \leq t < t_2$, r_t^h is the solution to $H_0 = N_t D_t^h(r_t^h)$, and for $t \geq t_2$, r_t^h solves $S - D^l(r_t^h - (1 - \beta)K) = N_t D_t^h(r_t^h)$. In each case, because $N_t = N_0 e^{g t + (\mu_t - 1)x}$, the resulting r_t^h is a differentiable function of x for any value of μ_t and is the same at $x = 0$ for any value of μ_t . Therefore, $\tilde{\mathbb{E}}r_t^h - \mathbb{E}_i r_t^h = O(x)$ as $x \rightarrow 0$ for all i , and the same then holds for p_1^h .

We now return to showing that if construction fails to occur at $t = 1$, then $r_0^h < r_1^h$. Using (2.9), we write $D_0^h(r_0^h) = N_0 f_0(r_0^h)$, and using (B.1), we write $D_1^h(r_1^h) = N_1 f_1(r_1^h)$. Without construction at $t = 1$, we have $N_0 f_0(r_0^h) = N_1 f_1(r_1^h)$. Note from (2.9) and (B.1) that $f_0 = f_1 + O(x)$ as $x \rightarrow 0$; this fact follows because $\tilde{\mathbb{E}}p_1^h - \mathbb{E}_i p_1^h = O(x)$ as $x \rightarrow 0$ for all i . Using $N_1 = N_0 e^{g + (\mu_1 - 1)x}$, we can conclude that $e^{g + (\mu_1 - 1)x} f_1(r_1^h) = f_1(r_0^h) + O(x)$ as $x \rightarrow 0$.

Because $e^{8+(\mu_1-1)x} > 1$ as $x \rightarrow 0$ and f_1 is decreasing, there exists an open $I_3 \ni 0$ such that for $x \in I_3$, $r_1^h > r_1^0$. This inequality is what we needed to show to prove that construction occurs at time 1, which is all that remained to prove that construction always occurs. We set $I = I_1 \cap I_2 \cap I_3$.

All of that proved that for $t > 0$, the effect of the shock z_t on r_t^h results from differentiating the equation $S - D^l(r_t^h - (1 - \beta)K) = D_t^h(r_t^h)$ with respect to x at $x = 0$. Doing so yields $-(D^l)'dr_t^h/dx = \mu_t D_t^h + (D_t^h)'dr_t^h/dx$, from which it follows that $dr_t^h/dx = -\mu_t D_t^h / ((D^l)' + (D_t^h)') = \mu_t r_t^h / (\epsilon_t^S + \epsilon^D)$. Similarly, the partial effect of the shock on current rents r_0^h , holding beliefs constant and letting N_0 change, is $\partial r_0^h / \partial x = r_0^h / (\epsilon_0^S + \epsilon^D)$. Putting together this partial effect with the one in Proposition 5 yields

$$\frac{dp_0^h}{dx} = \frac{r_0^h}{\epsilon_0^S + \epsilon^D} + \sum_{t=1}^{\infty} \left(\frac{\epsilon_0^S + (1 - \chi)\epsilon^D}{\epsilon_0^S + \epsilon^D} \tilde{\mu}_t + \frac{\chi\epsilon^S}{\epsilon_0^S + \epsilon^S} \bar{\mu}_t \right) \frac{\beta^t r_t^h}{\epsilon_t^S + \epsilon^S},$$

where $\tilde{\mu}_t$ is the most optimistic belief of μ_t and $\bar{\mu}_t$ is the average belief of μ_t . Because all residents agree that $\mu_0 = 1$, we may rewrite this expression as

$$\frac{dp_0^h}{dx} = \sum_{t=0}^{\infty} \left(\frac{\epsilon_0^S + (1 - \chi)\epsilon^D}{\epsilon_0^S + \epsilon^D} \tilde{\mu}_t + \frac{\chi\epsilon^S}{\epsilon_0^S + \epsilon^S} \bar{\mu}_t \right) \frac{\beta^t r_t^h}{\epsilon_t^S + \epsilon^S}.$$

The text defines the mean persistence of a persistence vector μ' to be $\mu' = \sum_{t=0}^{\infty} \mu_t \beta^t r_t^h (\epsilon_t^S + \epsilon^D)^{-1} / \sum_{t=0}^{\infty} \beta^t r_t^h (\epsilon_t^S + \epsilon^D)^{-1}$. We use this definition, and divide through by $p_0 = \sum_{i=0}^{\infty} \beta^i r_i^h$, which holds at $x = 0$, to derive

$$\begin{aligned} \frac{d \log p_0^h}{dx} &= \left(\sum_{t=0}^{\infty} \beta^t r_t^h \right)^{-1} \sum_{t=0}^{\infty} \left(\frac{\epsilon_0^S + (1 - \chi)\epsilon^D}{\epsilon_0^S + \epsilon^D} \tilde{\mu} + \frac{\chi\epsilon^S}{\epsilon_0^S + \epsilon^S} \bar{\mu} \right) \frac{\beta^t r_t^h}{\epsilon_t^S + \epsilon^S} \\ &= \left(\frac{\epsilon_0^S + (1 - \chi)\epsilon^D}{\epsilon_0^S + \epsilon^D} \tilde{\mu} + \frac{\chi\epsilon^S}{\epsilon_0^S + \epsilon^S} \bar{\mu} \right) \frac{1}{\tilde{\epsilon}^S + \epsilon^D}, \end{aligned}$$

where we have used the definition of the long-run supply elasticity $\tilde{\epsilon}^S$ given in the text. This equation for $d \log p_0^h / dx$ matches (2.12) in Proposition 6.

B.2.4 Proof of Implication 7

We demonstrate a limiting case in which $\epsilon_0^S = \infty$ while $\tilde{\epsilon}^S < \infty$. Let $D^l(r) = br^{-\epsilon^l}$ for some constant $b > 0$. Consider the limit as $b \rightarrow 0$. We know that $r_t^h \geq (1 - \beta)K$ because $r_t^h = r_t^l + (1 - \beta)K$ and $r_t^l \geq 0$. Define N^* to be the value of N_t that solves the equation $S = D_t^h((1 - \beta)K)$, where D_t^h is given by (B.1). For $N_t < N^*$, housing demand fails to exceed available land at the minimum rent, and there is no demand for land in the limit, so the market clearing rent must be $r_t^h = (1 - \beta)K$ while $H_t < S$. By (2.4), $\epsilon_t^S = \infty$ in this case. But for $N_t > N^*$, demand exceeds supply at the minimum rent, so $r_t^h > (1 - \beta)K$ and $H_t > 0$, leading to a finite elasticity. Since N_t grows at a constant rate g , for any $N_t < N^*$ we have $\epsilon_0^S = \infty$ but $\tilde{\epsilon}^S < \infty$.

B.2.5 Proof of Implication 8

Disagreement amplification Δ equals

$$\Delta = \frac{\epsilon_0^S + (1 - \chi)\epsilon^D}{\epsilon_0^S + \epsilon^D} \frac{\tilde{\mu} - \bar{\mu}}{\tilde{\epsilon}^S + \epsilon^D}.$$

We calculate this difference from subtracting from (2.12) the counterfactual in which we substitute $\bar{\mu}$ for $\tilde{\mu}$. Define $\bar{N}_0^*(\chi)$ to be the value of development (which determines the supply elasticities; see above) that maximizes Δ . When $\chi = 1$, Δ is 0 in the limits as $\bar{N}_0 \rightarrow 0$ and $\bar{N}_0 \rightarrow \infty$, because $\tilde{\epsilon}^S = 0$ in the first case and $\epsilon_0^S = 0$ in the second. But $\Delta > 0$ for $\chi = 1$, so $0 < \bar{N}_0^*(1) < \infty$ by continuity. But $\bar{N}_0^*(\chi)$ is continuous in χ as long as it exists and is finite, so there must exist $\chi^* < 1$ such that for $\chi^* \leq \chi \leq 1$, $\bar{N}_0^*(\chi)$ exists and is finite.

B.2.6 Proof of Implication 9

When $\chi = 1$, the limit as $\bar{N}_0 \rightarrow \infty$ of $d \log p_0^h / dx$ is $\bar{\mu} / \epsilon^D$. For any $0 < \bar{N}_0 < \infty$, we can choose $\tilde{\mu}$ to be large enough so that the price change given by (2.12) is larger than $\bar{\mu} / \epsilon^D$, because this price change becomes arbitrarily large with $\tilde{\mu}$. By continuity, we can do the same for some $\chi < 1$.

B.3 Construction equation

By the definition of supply elasticity, the change in the log housing stock is $\epsilon_0^S d \log r_0^h / dx$. The total effect of the shock on rents combines the effect in the end of the proof of Proposition 5 and the direct effect of the shock on N_0 derived in the proof of Proposition 6. It is $dr_0^h / dx = r_0^h / (\epsilon_0^S + \epsilon^D) - \chi \epsilon^D (\partial \beta \tilde{E} p_1^h / \partial x - \partial \beta \bar{E} p_1^h / \partial x) / (\epsilon_0^S + \epsilon^D)$. We substitute in for the beliefs from the Proof of Proposition 6 and divide through by r_0^h , and then multiply by ϵ_0^S to get

$$\frac{d \log H_0}{dx} = \frac{\epsilon_0^S}{\epsilon_0^S + \epsilon^D} \left(1 - \frac{\chi \epsilon^D}{\tilde{\epsilon}^S + \epsilon^D} \rho (\tilde{\mu} - \bar{\mu}) \right), \quad (\text{B.2})$$

where $\rho \equiv p_0^h / r_0^h$ is the price-rent ratio of the city before the shock at $x = 0$.

Appendix C

Appendix to Chapter 3

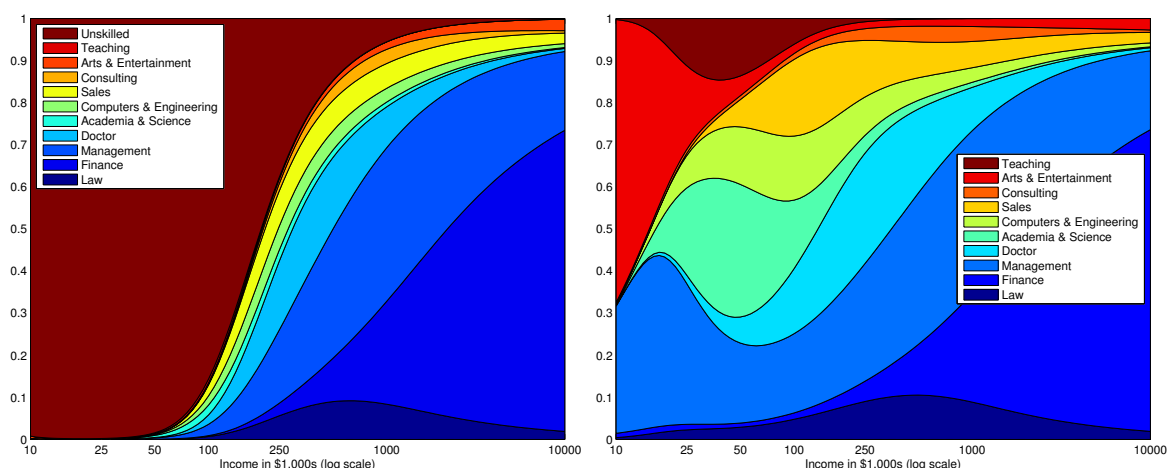
C.1 Alternative Finance Calibration

The one profession for which our approach to estimation worked poorly was Finance because of the large discrepancies between the Harvard and IRS data. In this appendix we maintain the same estimates of shares in different professions of the previous appendix, but fit the Finance-conditional distribution of income to fit only the Harvard data and disregard the IRS data as a robustness check on our conclusions.

Figure C.1 shows the results of this alternative calibration. The left shows the conditional profession distribution including the “Unskilled” category. An important thing to note is that the fraction of individuals captured by our professions falls at lower incomes. For example, at \$250k we now only account for about 70%. However, it is now monotonically increasing. The reason is that our alternative estimation lends Finance a much thicker upper tail (as this is what fits the Harvard data) and thus shifts many financiers from modestly wealthy income levels where there are a reasonable number of individuals in our original estimation to very high incomes that were much thinner.

The effects of this change on the professional distribution conditional on being among our professions is shown more clearly in the right panel of the figure. There we see that the new estimation makes only a modest difference in depressing Finance’s share at the

Figure C.1: *Income Distributions Under Alternative Finance Calibrations*



Notes: Conditional profession distribution by income for skilled individuals with (left) and without (right) the inclusion of “unskilled” professions under the alternative, Harvard-matching estimation of the Finance income distribution.

“thicker” incomes below \$1m. However, Finance’s share dramatically blows up at higher incomes, coming to dominate incomes above about \$2m overwhelmingly.

Figure C.2 shows the impact this has on ATEM and MTEM optimal taxes. The results are very close to those in the text on which we focus qualitatively, except that now under both literature calibration rates rise monotonically at higher incomes and there is a much smaller difference between the pro mgmt calibration and the anti mgmt calibration. Rates at very high incomes are somewhat higher under all calibrations other than Tea Party.

We omit the reference distributions by ability levels, as these are essentially the same except that Finance is now more unequal and thus at very high incomes overtakes Management as the most lucrative occupation. Optimal marginal tax rates under our baseline structural model (as in Subsection 3.5.2) also hardly change, except that pro mgmt comes to resemble anti mgmt even more closely. More revealing is the difference in the impact of externalities as compared to intensive margin labor supply elasticities on the structure and level of optimal marginal tax rates, as in Subsection 3.5.3. This analysis is shown in Figure C.3. Compared to Figure 3.8, our results here show that rates are much more clearly

Figure C.2: ATEM and MTEM Policies Under Alternate Finance Calibration

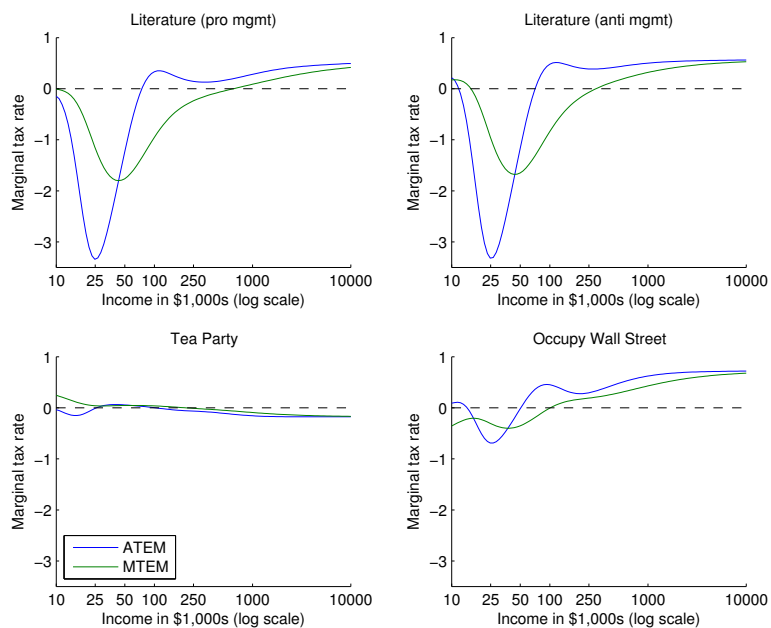
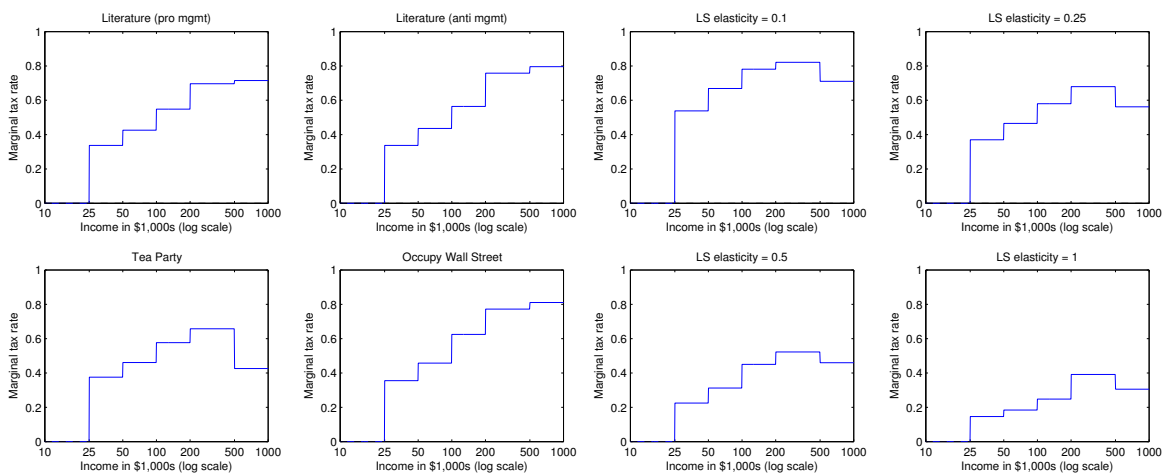


Figure C.3: “Horseshoe” Under Alternate Finance Calibration



progressive under all but the Tea Party calibrations and much more extremely regressive under the Tea Party calibration. Overall levels of rates move only slightly, both for the externality-varying calibrations and the elasticity-varying calibrations. However, now, under the elasticity-varying calibrations the extent of regressivity is now more-or-less constant across values of the elasticity. This alternative calibration thus makes even stronger our claim that externalities have a stronger effect on the structure of optimal taxes than does the intensive elasticity of labor supply.

Our results on quantitative welfare gains (as in Subsection 3.5.4) are quite similar, except that they are somewhat (roughly 25%) larger, especially under the Occupy calibration. More interesting is the impact of the new calibration on the effect of the Reagan tax reforms. All impacts on the allocation of talent are in the same direction but a bit (10-20%) larger. The impacts of the reforms on social welfare and GDP are more negative, and more consistent between the two literature calibrations. However the largest change is on the impact of the reforms on inequality. This *doubles* as a result of the new calibration, so that now the Reagan reforms account for 42% of the 1979-1993 and 30% of the 1979-1997 change in the share of income earned by talented individuals in the top 1%.

C.2 Externality Share Calibration

In this appendix we discuss in greater detail our strategy for calibrating the externality shares of various professions based on the literature.

C.2.1 Law and Computers/Engineering

Murphy *et al.* (1991) calculate the direct spillovers of both lawyers and engineers using cross-country regressions. We use their preferred estimates that are restricted to the 55 countries in which more than 10,000 students are in college. Their result is that a 1% point increase in the share of students studying law lowers real per capita GDP growth rates by 0.078% points, and that a 1% point increase in the share of students studying engineering increases real per capita GDP growth rates by 0.054% points.

We convert these estimates to per-dollar externalities using our estimates of the income distributions in law and engineering, as well as OECD data on the share of students studying each field in the United States. We interpret the spillovers of these professions on GDP growth as a one-time, static reduction (or improvement) to GDP. This interpretation is consistent with our static model that does not consider economic growth.¹ All of the following data comes from 2005.

The OECD data tell us that 2.4% of students study law (51,000 out of 2,100,000). US GDP in 2005 was \$12.4 trillion. The total externality from lawyers is therefore

$$\underbrace{-0.078}_{\text{MSV estimate}} \cdot 0.024 \cdot \$12.4 \text{ trillion} = -\$23 \text{ billion.}$$

To compute the externality share e_p for law, we divide this figure by the aggregate income for law. We focus just on the aggregate income of the “talented” individuals for which we calibrated income distributions in Section 3.3.1. We attribute all of the spillovers calculated by MSV to this group of lawyers. According to our calibrated income distribution, the mean income of these lawyers is \$337,000 and there are 325,000 of them. Their total income is therefore \$110 billion, and the per-dollar externality from law is then

$$e_{\text{Law}} = -\frac{\$23 \text{ billion}}{\$110 \text{ billion}} = -0.21.$$

Analogous calculations yield our estimate for the externality share of engineering. The OECD data tell us that 10.7% of students study engineering (230,000 out of 2,100,000). The total externality from engineers is therefore

$$\underbrace{0.054}_{\text{MSV estimate}} \cdot 0.107 \cdot \$12.4 \text{ trillion} = \$72 \text{ billion.}$$

We again attribute all of the spillovers calculated by MSV to the “talented” engineers whose income distributions we estimated in Section 3.1. According to that distribution, the mean

¹An alternate approach is to interpret the spillovers as permanent shocks to GDP, which affect welfare in current and all future years. We would then compute the present value of these spillovers. This methodology is more appropriate in many ways but our static model is ill-suited to analyze it.

income of these engineers is \$172,000 and there are 560,000 of them. Their total income is therefore \$96 billion, and the per-dollar externality from engineering is then

$$e_{\text{Engineering}} = \frac{\$72 \text{ billion}}{\$96 \text{ billion}} = 0.75.$$

C.2.2 Management

The externality share for management largely determines top tax rates because management is the most heavily represented occupation at top tax rates. Two strands in the literature offer competing views on the externalities of management. We calculate optimal tax rates using management externalities calibrated from each strand of the literature. Our paper offers a framework for analyzing how the conclusions from that literature quantitatively affect optimal tax rates.

The first half of the literature argues that executive compensation reflects a good deal of rent seeking. According to this story, Chief Executive Officer (CEO) compensation shifts resources from shareholders to managers in ways that do not actually reflect the CEO's marginal product. Papers that try to document such behavior include Bertrand and Mullainathan (2001) and Malmendier and Tate (2009). Piketty *et al.* (2014) argue that 60% of the CEO earnings elasticity with respect to taxes represents this rent seeking behavior. The other 40% is genuine labor supply. This literature therefore gives the estimate

$$e_{\text{Management}} = -0.6.$$

The other half of the literature argues that the increases in CEO compensation over the last 30 years are due not to rent seeking but to fundamental labor market factors, such as increasing firm size (Gabaix and Landier, 2008). Furthermore, the compensation patterns highlighted by the aforementioned papers as evidence of rent seeking can be rationalized as efficient dynamic contracts. Edmans and Gabaix (2009) offer a survey of this literature. The conclusion from this perspective is that CEO pay reflects the marginal product of these

managers, and that therefore

$$e_{\text{Management}} = 0.$$

We assume all management has the same externality share as CEOs.

C.2.3 Academia/Science

We estimate that the aggregate income for talented workers in Academia and Science is \$67 billion (\$107,000 per worker times 626,000 workers). Aggregate income is therefore \$67 billion.²

Murphy and Topel (2006) estimate that the gains of medical research from 1970 to 2000 were \$3.2 trillion annually. This spillover is likely the largest externality from academia and science, but note that this estimate is still conservative in assuming no gains accrue from any activity in academia other than medical research. The resulting externality share for all of this occupation is then

$$e_{\text{Academia/Science}} = \frac{\$3.2 \text{ trillion}}{\$67 \text{ billion}} \approx 48.$$

An alternative benchmark for the value of scientific research is much narrower and focuses only on the spillovers of universities to profits made by geographically proximate firms, thus neglecting other profit and consumer welfare spillovers. This narrower measure is studied by Jaffe (1989). His model allows university research to have a direct effect on commercial patents as well as an indirect effect through influencing industrial R&D. His preferred estimate is that the total elasticity of patents with respect to university research is .6. The direct elasticity of industrial R&D on patents is .94 and industrial R&D expenditures are 6 times larger than university research expenditures. Therefore a marginal dollar in university research is equivalent to $\frac{6 \cdot 6}{.94} = 3.83$ dollars spent on industrial R&D in terms of resulting patents.

According to the National Science Foundation, \$45 billion was spent on university R&D

²According to <http://www.researchamerica.org/uploads/healthdollar12.pdf>, total spending on medical research was \$80 billion in 2001. So the prior number is the right order of magnitude.

in 2005. Using the estimate from Jaffe, we conclude that the total spillover from this activity was \$172.35 billion. Dividing this by our \$67 billion of salary yields a much lower 2.6 externality share.

The largest difference between these measures is how broadly they attempt to measure the spillovers from academic research. While we are more sympathetic to the Murphy and Topel accounting, using 48 as the externality share of Academia/Science would cause every other factor to be swamped by the need to get more individuals into this field. Perhaps this is appropriate, but we felt that, given the smaller Jaffe numbers, a reasonable compromise was to use 5, twice the Jaffe number but still an order of magnitude smaller than that from Murphy and Topel.

C.2.4 Consulting

A body of work has demonstrated that better management practices increase firm productivity. We interpret the marginal product of consulting as teaching firms these management practices. Some examples in which management practices increase productivity are in adopting new technologies (Bloom *et al.*, 2012), decreasing mortality in cardiac care centers (Bloom *et al.*, 2011; McConnell *et al.*, 2013), and teaching and research output in universities (McCormack *et al.*, Forthcoming).

Given these wide-ranging examples where management quality increases productivity, our prior is that consulting fully captures this marginal product with the fees they charge for their services. Some direct evidence on the relationship between consulting fees and increased productivity comes from Bloom *et al.* (2013). The authors randomly assign consultants to plants in India. The market fees for the consultants' time (which they did not actually charge) was on average \$300,000 per plant, and they increased the productivity of each plant on average by \$250,000. We interpret these numbers as in line with our prior from the larger field of research that consultants' fees match their marginal product, and thus set

$$e_{Consulting} = 0.$$

C.2.5 Teaching

Chetty *et al.* (2013a,b) measure the extent to which higher quality elementary and middle school teachers raise the eventual earnings of their students. Their main result is that a 1 standard deviation increase in teacher quality raises eventual student earnings by 1.34%. We use their results to estimate the surplus social benefit of the entry of talented individuals into the teaching profession.

Chetty *et al.* measure teacher quality with student test scores. For our purposes, we thus need an estimate of the extent to which talented workers would raise student test scores. Teach For America is a program in which college graduates from selective colleges teach for one or two years. Clark *et al.* (2013) estimate that Teach For America teachers raise math test scores by 0.07 standard deviations. They compare 136 Teach For America teachers to teachers teaching in similar classrooms at the same schools.

Most of the teachers in the Clark *et al.* (2013) study are in middle schools, so we use the middle school numbers from Chetty *et al.*. Chetty *et al.* (2013a) find that 1 standard deviation in middle school math teaching quality corresponds to a 0.134 standard deviation increase in student test scores (Table 2 of that paper). We conclude that talented workers (represented here by TFA members) are $0.52 = 0.07/0.134$ standard deviations of teacher quality better than the average teacher.³

Chetty *et al.* (2013b) calculate that the present value of future earnings for a middle school student is \$468,000 in 2005 dollars.⁴ They also calculate that the average number of students in a class is 28.2. The surplus social value of the entry of a talented worker into teaching is

$$0.0134 \cdot 0.52 \cdot 28.2 \cdot \$468,000 = \$92,000.$$

³We are assuming that the average treatment effect on student earnings of a teacher whose students have higher test scores is the same for all teachers as it is for Teach For America teachers. It is possible that TFA teachers raise test scores via a different channel than other teachers, in which case the treatment affect would not be the same. Similar assumptions are made in extrapolation to earnings.

⁴They use a 5% discount rate and assume that earnings grow 2% annually.

Our calibrated income distribution for “talented” teachers implies that the mean income for such teachers is \$83,000. According to Chetty *et al.* (2013b), the average income for teachers is \$50,000. We interpret this figure to represent the marginal social product of the average teacher. The average teacher’s externality share is 0 under this assumption. A talented teacher’s total social product is therefore $\$50,000 + \$92,000 = \$142,000$, and the resulting externality share is

$$e_{Teaching} = \frac{\$142,000}{\$83,000} - 1 = 0.71.$$

C.2.6 Arts/Entertainment

Several quantitative studies (Florida and Mellander, 2010; Hellmanzik, 2010) and much qualitative work (Jacobs, 1961) indicates that artists have important spillovers to the value of property in urban areas. It seems plausible that cultural production has other positive externalities such as social cohesion (Lazear, 1999) and un-captured consumer surplus in a manner similar to technological innovations. However these effects have proven much harder to measure (Hausman, 2012) and we saw no way to use those that have been done to generate a number. Furthermore, unless this number is quite large it has little affect on our results, except to somewhat dampen the subsidy on entering the middle class in Subsection 3.3.3, because so few individuals enter Arts/Entertainment. We thus reluctantly used a 0 externality share.

C.3 General Ability Model

First we discuss an analytical result for the case when $N = 2$ under the general ability assumption, but without further restrictions on functional forms for distributions.

Proposition 9. *Suppose that G_1 first-order stochastically dominates G_2 . If $e_1 < e_2$ then beginning from ATEM, there is a first-order welfare gain from raising the marginal tax rate at each wage level w in the support of the wage distribution for which $\partial\Theta(w)$ includes individuals switching either from*

or to a wage level including individuals from both professions. On the other hand, if $e_1 > e_2$ then beginning from ATEM, there is a strictly positive first-order welfare gain from lowering the marginal tax rate at each wage level w in the support of the wage distribution for which $\partial\Theta(w)$ includes individuals switching either from or to a wage level including individuals from both professions.

Proof. Suppose that $T(w) = T_0 - E[e_{p^*}|\Theta(w)]w$. Consider any w meeting the description in the proposition statement. Partition $\partial\Theta(w)$ into two sets: $\widetilde{\partial\Theta}(w)$ and $\overline{\partial\Theta}(w)$ where $\widetilde{\partial\Theta}(w)$ includes only types switching between wage levels at least one of which includes a positive density of individuals of each profession and $\overline{\partial\Theta}(w)$ includes only individuals switching between wage levels occupied each by a sole profession. By assumption, $\widetilde{\partial\Theta}(w)$ is non-empty; let

$$\lambda(w) \equiv \frac{\int_{\widetilde{\partial\Theta}(w)} f(\theta) d\theta}{\int_{\widetilde{\partial\Theta}(w)} f(\theta) d\theta + \int_{\overline{\partial\Theta}(w)} f(\theta) d\theta}.$$

Then

$$\begin{aligned} E[\Delta T + \Delta E|\partial\Theta(w)] &= \lambda(w)E\left[T(w_{ip}) + e_{ip}w_{ip} - T(w_{iq}) - e_{iq}w_{iq}|\widetilde{\partial\Theta}(w)\right] + \\ &\quad [1 - \lambda(w)]E\left[T(w_{ip}) + e_{ip}w_{ip} - T(w_{iq}) - e_{iq}w_{iq}|\overline{\partial\Theta}(w)\right]. \end{aligned}$$

Note that the second term is 0 by the argument in the proof of Proposition 8. On the first term, note that for any type in $\widetilde{\partial\Theta}(w)$ we have

$$\begin{aligned} T(w_{ip}) + e_{ip}w_{ip} - T(w_{iq}) - e_{iq}w_{iq} &= \\ -E[e_{p^*}|\Theta(w_{ip})]w_{ip} + e_{ip}w_{ip} + E[e_{q^*}|\Theta(w_{iq})]w_{iq} - e_{iq}w_{iq} &< (>)0 \end{aligned}$$

if $e_1 < (>)e_2$ as $p = 1$ and $q = 2$ always because no individual has a higher income in profession 2 than 1 by the general income assumption and first-order stochastic dominance. \square

While this result does not directly imply that marginal rates are higher (lower) at every point when the high-paying profession has more negative (positive) externalities than would be the case under Proposition 8, they are strongly suggestive in this direction. Furthermore the result applies only when there are two professions. Thus to verify whether in a realistic

model with many professions externalities do end up making a greater impact on optimal tax rates than under ATEM requires an empirical calibration.

In particular, we adapt our general model to a logit discrete choice framework to describe agents' career selection conditional on ability level. There are N professions, indexed $p = 1, \dots, N$, and a unit mass of individuals, indexed by i . Each individual i is characterized by general ability $a_i \in (0, 1)$ and a vector of psychic incomes across professions, $\psi_i = \{\psi_{ip}\}_{p=1}^N$. Each profession has wage "reference distribution" $G_p(\cdot)$ mapping ability to marginal product. Utility is quasilinear in consumption and psychic income, and the marginal utility of consumption is normalized to one. In our basic model, labor supply is fixed at one, so that i 's income in profession p , denoted y_{ip} , is simply $G_p(a_i)$. The retention function $R(\cdot)$ maps pretax income to consumption.

Psychic income consists of two components, a profession-specific mean (independent of ability) $\hat{\psi}_p$, and an individual-specific residual ϵ_{ip} . Individual i selects the profession $p^*(i)$ that maximizes utility:

$$\begin{aligned} p^*(i) &= \arg \max_p \{R(y_{ip}) + \psi_{ip}\} \\ &= \arg \max_p \{R(G_p(a_i)) + \hat{\psi}_p + \epsilon_{ip}\}. \end{aligned} \quad (\text{C.1})$$

(Under quasilinear utility, externalities imposed by other agents do not change the relative attractiveness of professions and so can be ignored.) Let $\epsilon_{ip} = \hat{\epsilon}_{ip}\beta(a_i)$, where $\hat{\epsilon}_{ip}$ is iid across professions, drawn from a standard Gumbel (extreme value) distribution, and $\beta(a_i)$ scales $\hat{\epsilon}_{ip}$ so that it is expressible in dollars (and thus is additive in utility). Then Equation C.1 can be written as

$$p^*(i) = \arg \max_p \left\{ \beta(a_i)^{-1} (R(G_p(a_i)) + \hat{\psi}_p) + \hat{\epsilon}_{ip} \right\}. \quad (\text{C.2})$$

The logit model implies that the probability that i selects a particular profession p' has a closed form expression:

$$\Pr(p' = p^*(i)) = \frac{\exp [\beta(a_i)^{-1} (R(G_{p'}(a_i)) + \hat{\psi}_{p'})]}{\sum_{p=1}^N \exp [\beta(a_i)^{-1} (R(G_p(a_i)) + \hat{\psi}_p)]}. \quad (\text{C.3})$$

Equivalently, this expression gives the share of agents with i 's ability level who select profession p' . If psychic incomes $\hat{\psi} = \{\hat{\psi}_p\}_{p=1}^N$ and reference distributions $G_p(\cdot)$ are known, we can calculate the share of agents at each ability level selecting each profession given any counterfactual wages; in particular, we can calculate the counterfactual distribution of professions given any tax function. The function $\beta(a)$ controls the sensitivity of agents with ability a to differences in incomes—higher $\beta(a)$ indicates less willingness to switch careers in the face of a given income spread.

C.4 Estimation

To estimate the model, we must estimate the function $\beta(a)$, the vector $\hat{\psi}$, and the reference distributions $G_p(\cdot)$. To make this task tractable, we make a key parametric assumption: we assume the standard deviation of utility from profession p across agents with ability a , which is $\beta(a)\sqrt{\pi/6}$, is proportional to the average marginal product across professions at that ability level.⁵ That is, $\beta(a) = \beta \cdot \mathbb{E}[G_p(a)]$, for some constant β . This assumption implies that lower earners are more sensitive to income differences across professions than high earners. In practice our results are not highly sensitive to this assumption—even using a constant value for $\beta(a)$ yields qualitatively similar results. We further assume that all professions have positive support across the ability distribution, and that the reference distributions $G_p(\cdot)$ have range $(0, \infty)$ for all professions.

Estimating β

To estimate β , we assume the distribution of ability among Harvard graduates is sufficiently narrow to be considered constant, denoted a_H . We estimate $\beta(a_H)$ using data about Harvard

⁵The standard Gumbel distribution has a variance of $\pi/6$.

graduates and changes in financial salaries, from which we can back out the value of the constant β .

The derivative of (C.3) with respect to $R(y_{fin}(a_H))$, simplified, is

$$\frac{\partial s_p(a_H)}{\partial R(y_{fin}(a_H))} = \beta(a_H)^{-1} (s_p(a_H) - s_p(a_H)^2), \quad (\text{C.4})$$

implying that we can calculate $\beta(a_H)$ using three figures: the change in post-tax salaries in finance, relative to other professions, over a given period ($\partial R(y_{fin}(a))$), the change in the share of Harvard graduates entering finance over the same period ($\partial s_p(a_H)$), and the mean share of Harvard graduates entering finance during this period ($s_p(a_H)$).

According to Goldin *et al.* (2013), from 1970 to 1990, the share of male Harvard graduates entering finance rose from 5% to 16%, so we use $\partial R(y_{fin}(a)) = 0.11$ and $s_p(a_H) = 0.105$.

To compute $\partial R(y_{fin}(a))$, we use Philippon and Reshef (2012) time series estimates of finance wages. They present several data series, which control for education and alternate occupational choice in different ways. We use the wage in finance relative to other industries (their Figure 1) for our calculations—the approximate values are listed in the first column of Table C.1 below. From Goldin *et al.* (2013), we know Harvard alumni in finance who graduated in 1988–1992 had an annual salary in 2005 of approximately \$615,000. Using Philippon and Reshef’s figures as an approximation for the real change in financial salaries, we compute real income in finance in the second column of Table C.1. In the following columns we convert these real values to nominal values, and use NBER’s TAXSIM to compute nominal post-tax income in each year, then convert back to real figures in the final column.⁶ This allows us to explore how the salary premium in finance varied between 1970 and 1990. If graduates in those years looked to salaries at the time of graduation when selecting careers, then they would have seen an increase from about \$225,000 to \$300,000, suggesting a value for $\partial R(y_{fin}(a))$ of \$75,000.

If, on the other hand, graduates correctly predicted what salaries would be, say, 15 years

⁶This procedure is clearly imperfect, as we do not have figures for the absolute change in financial salaries, which would be appropriate for computing tax burdens. We use a conservative estimate and a wide range for robustness checks to ensure our results are not sensitive to this imperfection.

Table C.1: *Real net income in finance over time*

	P&R Fig 1	Real Income	Nom. Income	Tax Burden	Nom. Consump	Real Consump
1970	1.05	379,853	75,971	31,260	44,711	223,553
1980	1	361,765	151,941	67,824	84,117	200,279
1985	1.1	397,941	218,868	94,634	124,234	225,879
1990	1.2	434,118	290,859	87,763	203,096	303,128
2000	1.6	578,824	509,365	195,686	313,679	356,453
2005	1.7	615,000	615,000	214,808	400,192	400,192

after they entered finance, the increase would instead be \$225,000 to \$400,000, suggesting $\partial R(y_{fin}(a)) = \$175,000$. Predicting wages only 10 years later would instead correspond to an increase of approximately \$150,000. We select a benchmark value for $\partial R(y_{fin}(a))$ of \$150,000, and we perform robustness checks using both \$75k and \$300k (see Appendix C.6 below). Using Equation C.4, we compute corresponding values of $\beta(a_H)$, and dividing by the mean salary among Harvard graduates across all professions (\$189,000 in the benchmark case of zero labor supply elasticity) we obtain the constant β .

Estimating psychic incomes $\hat{\psi}$ and reference distributions $G_{p'}(\cdot)$

Psychic incomes and reference distributions are solved jointly, using a computation that is most easily understood using a simple example with $N = 3$ professions, $\{a, b, c\}$, and $K = 2$ discrete levels of income, $\{y_\ell, y_h\}$. In this simplified example, observable data would consist of the population shares in each profession at each level of income,

$$\begin{pmatrix} f_{a\ell} & f_{b\ell} & f_{c\ell} \\ f_{ah} & f_{bh} & f_{ch} \end{pmatrix}, \quad (\text{C.5})$$

representing five degrees of freedom (since the shares must sum to 100%), or, in general, $NK - 1$ degrees of freedom. Letting a be the reference profession, the vector of psychic incomes in this case is $\hat{\psi} = (\psi_b \ \psi_c)$. In general, this vector will have length $K - 1$. The reference distributions $G_{p'}(\cdot)$ in this case are approximated by a 1×3 matrix $\hat{G} = (g_a \ g_b \ g_c)$, where for example g_a represents the ability percentile below which an individual in profession a earns y_ℓ , and above which an individual earns y_h . More generally, given

K discrete income levels, the matrix characterizing the reference distributions will be of dimension $(K - 1) \times N$, representing the $K - 1$ percentiles partitioning the ability space into income levels within in each profession. Note that the total number of free parameters between $\hat{\psi}$ and \hat{G} is $(N - 1) + (K - 1) \times N = NK - 1$, matching the number of independent points in the data.

The system can be solved by constructing a function $s(\hat{\psi}, \hat{G})$ to compute the shares \hat{f} that would arise for given values of $\hat{\psi}$ and \hat{G} , and then searching for the $\hat{\psi}$ and \hat{G} such that these resulting shares match the data observed in (C.5). This function can be constructed as follows. Place the elements of \hat{G} in a vector in increasing order, thereby partitioning the ability space into $(K - 1) \times N + 1$ groups, each specifying a vector of salaries faced across professions. For example, if a candidate \hat{G} has $g_b < g_a < g_c$, the vector

$$\begin{pmatrix} g_b \\ g_a \\ g_c \end{pmatrix}$$

partitions the ability distribution into four groups. Individuals with ability below g_b face y_ℓ in each profession. Individuals with ability between g_b and g_a would earn y_ℓ in professions a and c , but y_h in profession b , and so on. This vector can be used to construct a $K \times N$ matrix in which element (m, n) represents the salary a member of group $m = 1 \dots 4$ would earn in profession n :

$$\begin{pmatrix} y_\ell & y_\ell & y_\ell \\ y_\ell & y_h & y_\ell \\ y_h & y_h & y_\ell \\ y_h & y_h & y_h \end{pmatrix}. \quad (\text{C.6})$$

This matrix, along with the candidate vector of psychic incomes $\hat{\psi}$ and an assumed retention function combine to characterize a matrix of utilities that an individual in each ability

partition would realize from each profession, except for the idiosyncratic component ϵ_{ip} :

$$\begin{pmatrix} u_{1a} & u_{1b} & u_{1c} \\ u_{2a} & u_{2b} & u_{2c} \\ u_{3a} & u_{3b} & u_{3c} \\ u_{4a} & u_{4b} & u_{4c} \end{pmatrix} = \begin{pmatrix} R(y_\ell) & R(y_\ell) + \psi_b & R(y_\ell) + \psi_c \\ R(y_\ell) & R(y_h) + \psi_b & R(y_\ell) + \psi_c \\ R(y_h) & R(y_h) + \psi_b & R(y_\ell) + \psi_c \\ R(y_h) & R(y_h) + \psi_b & R(y_h) + \psi_c \end{pmatrix}. \quad (\text{C.7})$$

We have assumed that $\beta(a')$ is proportional to the average marginal product of an individual with ability a' across professions, so that, for example, $\beta_2 = \beta(y_\ell + y_h + y_\ell)/3$. We can thus compute the logit shares of each partition that would select each profession. For example the share of partition 2 selecting profession b is given by

$$s_{2b} = \frac{\exp[(R(y_h) + \psi_b)/\beta_2]}{\exp[(R(y_\ell))/\beta_2] + \exp[(R(y_h) + \psi_b)/\beta_2] + \exp[(R(y_\ell) + \psi_c)/\beta_2]}.$$

The population share in partition 2 is $g_a - g_b$, so individuals in partition 2 selecting profession b make up fraction $s_{2b}(g_a - g_b)$ of the total population. Since these individuals earn y_h , as specified in (C.6), they contribute to \hat{f}_{bh} . Members of partitions 3 and 4 also contribute to \hat{f}_{bh} , therefore the total implied share of y_h -earners in profession b is

$$\hat{f}_{bh} = s_{2b}(g_a - g_b) + s_{3b}(g_c - g_a) + s_{4b}(1 - g_c).$$

The full matrix of implied \hat{f} shares is computed by weighting the logit shares arising from (C.7) by the population fraction in each ability partition, then the resulting shares are aggregated using the income levels in (C.6) as indices. To perform the numerical computation, we designate Law as the reference profession, and we set $K = 50$. For numerical simplicity, we use a flat tax approximating the US income tax for $R(\cdot)$ (results are not sensitive to this assumption). We use the numerical solver Knitro (Byrd *et al.*, 2006) to solve the system, dropping one element of the shares matrix to account for the $NK - 1$ degrees of freedom. In practice we search over the *spread* between thresholds in each profession, so that the search variables can be constrained to be nonnegative, and we impose the linear inequality constraint that these spreads sum to less than one. We also provide the solver with analytic gradients. Code for this optimization is available upon request.

C.5 Solving for the optimal tax function

To find the optimal tax function, we must specify a social welfare function given the model above. As described above, i 's utility is given by

$$\max_p \{R(G_p(a_i)) + \hat{\psi}_p + \hat{\epsilon}_{ip}\beta(a_i)\}.$$

Therefore a utilitarian social planner will select $R(\cdot)$, or equivalently, $T(\cdot)$, to maximize the expectation of this expression across all agents. To do this, it is sufficient to calculate the expected utility of all agents with a given ability a' , across all professions:

$$\begin{aligned} \mathbb{E}_{i,p} \left[\max_p \{R(G_p(a_i)) + \hat{\psi}_p + \hat{\epsilon}_{ip}\beta(a_i)\} \mid a_i = a' \right] = \\ \beta(a') \mathbb{E}_{i,p} \left[\max_p \{ \beta(a')^{-1}(R(G_p(a')) + \hat{\psi}_p) + \hat{\epsilon}_{ip} \} \right]. \end{aligned} \quad (\text{C.8})$$

We cannot estimate (C.8) directly, as we cannot observe $\hat{\epsilon}_{ip}$. However given our estimated model, we can calculate $\beta(a')^{-1}(R(G_p(a')) + \hat{\psi}_p)$ in a given profession, and by assumption we know the distribution of $\hat{\epsilon}_{ip}$. Williams (9) and Small and Rosen (1981) demonstrate that the iid extreme value distribution of $\hat{\epsilon}_{ip}$ implies (C.8) is equal to

$$\beta(a') \ln \left(\sum_{p=1}^N \exp \left(\beta(a')^{-1}(R(G_p(a')) + \hat{\psi}_p) \right) \right) + C, \quad (\text{C.9})$$

where C is a constant. Since we are interested in the differences in (C.9) under different tax functions, the constant C can be ignored.

Adding a redistributive motive

The above approach assumes all agents have the same marginal utility of consumption, and thus it is sufficient to evaluate tax functions in dollar terms. To incorporate a redistributive motive, we assume the marginal utility of consumption is lower for agents with higher

ability (who have higher average utility). Equivalently, we will place a lower welfare weight on consumption for high ability agents; we will use $\alpha(a)$ to denote the welfare weight assigned to ability a . To preserve the structure of the logit discrete choice framework, we continue to assume agents of each ability level have constant marginal utility of consumption. Nevertheless, we will allow the planner's welfare weights to depend on expected levels for each ability at the optimum. Letting $\bar{U}(a)$ denote the average utility at ability a , we use $\alpha(a) = \bar{U}(a)^{-\gamma}$, so that $\gamma \geq 0$ governs the curvature of the social welfare function, with the special case $\gamma = 0$ corresponding to no redistributive taste, and $\gamma = .95$ taken to approximate logarithmic preferences over agents' utility (our baseline redistributive case). We use this instead of logarithmic preferences (corresponding to $\gamma = 1$) as it makes little substantive difference and speeds computational convergence.

Adding a labor-leisure choice

So far we have assumed labor supply is inelastic with respect to the tax rate; we now relax that assumption. Whereas we previously assumed utility was given by $R(y_{ip}) + \psi_{ip}$, we will now suppose utility is given by

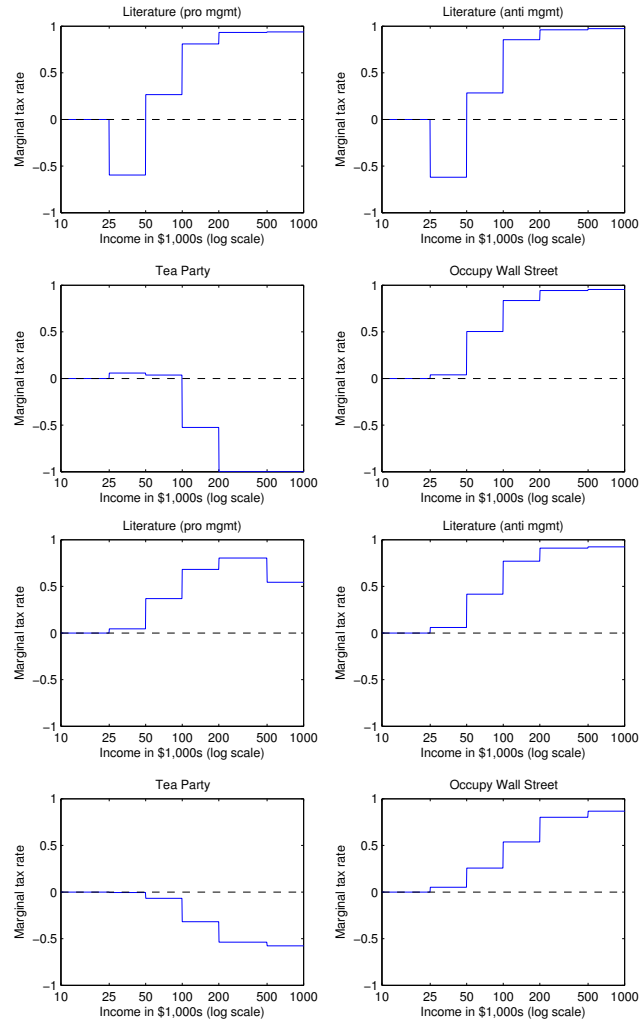
$$R(y_{ip}) - \frac{h^{1/\sigma+1}}{1/\sigma+1} + \psi_{ip},$$

where the agent now selects the labor supply h that maximizes utility. Pre-tax earnings y_{ip} is equal to the product of i 's supplied labor and i 's private marginal product in this profession, w_{ip} . The optimal labor supply choice solves $h^{1/\sigma} = R'(y)w_{ip} = (1 - \tau)w_{ip}$, where τ is the marginal tax rate on earnings y . Optimal labor supply is equal to

$$h_{ip} = (w_{ip}(1 - \tau))^\sigma$$

so σ is the elasticity of labor supply. We can use the equation to infer marginal product w_{ip} from observed earnings y_{ip} . As previously noted, $\beta(a_i)$ is assumed to scale with the

Figure C.4: Optimal Tax Rates Under Different Elasticity Calculation



Notes: Computational optimal piece-wise-constant marginal tax rates under robustness check where the increase in finance salaries is assumed to be \$75,000 (left) or \$300,000 (right).

mean earnings of ability a_i across professions—since we do not want this parameter to be endogenous to the tax function, in the presence of a labor-leisure choice we assume $\beta(a_i)$ scales with the earnings an agent with ability a_i would earn given her labor choice in each profession under a *laissez-faire* tax regime.

C.6 Alternative Elasticity Value

Table C.2: *Welfare Gains Under Alternate Elasticity Value*

	Gains from 1st best in Social Welfare (GDP)	Gains from 2nd best	Share of potential from 2nd best
\$75k increase			
Literature (pro mgmt)	35% (47%)	2.7% (1.3%)	7.7% (2.8%)
Literature (anti mgmt)	40% (52%)	5.1% (3.8%)	12% (7.2%)
Tea Party	1.7% (4.6%)	.97% (3.5%)	56% (76%)
Occupy Wall Street	7.8% (9.8%)	2.9% (.95%)	37% (9.6%)
\$300k increase			
Literature (pro mgmt)	5.5% (14%)	.62% (-.61%)	11% (-4.3%)
Literature (anti mgmt)	7.3% (18%)	1.6% (.61%)	22% (3.4%)
Tea Party	.58% (3.1%)	.25% (2.4%)	44% (76%)
Occupy Wall Street	1.9% (2.6%)	.88% (-.38%)	45% (-14%)

Notes: Quantitative welfare gain from non-linear income taxation (2nd best), compared to optimal profession-specific taxation (first best), relative to laissez-faire tax regime, under robustness checks where the increase in finance salaries is assumed to be either \$75,000 or \$300,000.

As discussed in Appendix Section C.4 above, there is ambiguity about the appropriate value of β , driven by uncertainty regarding the accuracy of Harvard graduates' predictions about the salary premium that would exist in finance years after graduation. Therefore we rerun our analysis under the assumption that expected finance salaries for Harvard graduates increased by, alternatively, \$75,000 or \$300,000 between 1970 and 1990, rather than our benchmark assumption of \$150,000. Results of this robustness check have little impact on our results other than the quantitative welfare gains and the quantitative (but not qualitative) impact of the Reagan tax reforms. An example is shown in Figure C.4, which gives the optimal tax rates analogous to the baseline scenario in Subsection 3.5.2 using

for both alternatives. The pattern of tax rates remains similar, though higher switching elasticities tend to magnify subsidies to the middle class under the Literature specifications and make top marginal tax rates more extreme under all specifications.

The muted sensitivity of these patterns to large changes in the switching elasticity is unsurprising, because optimal Pigouvian taxation is unaffected by the size of elasticities; only the quantitative benefits of imposing such taxation are impacted. The new results for the quantitative welfare gains are shown in Table C.2. The welfare gains from optimal taxation under the high elasticity robustness check (\$75k increase) are roughly twice as large as our benchmark specification, while the gains under the low elasticity alternative are one third to one half as large. Thus higher elasticity values essentially scale up the importance of the results proportionally.

C.7 Allocation of Talent in the General Ability Model

In this appendix we discuss how talent is allocated under various regimes in our structural model developed in Section 3.5. Figure C.5 shows the allocation of talent in the absence of taxation. Sales and Academia/Science attract a large number of individuals at the bottom of the ability distribution. A variety of professions are represented throughout the mid-range of the income distribution. But at the top and especially the very top Finance and Consulting are dominant and the “middle class” professions of Doctor, Academia/Science and Computers/Engineering largely die off.

Matters are radically different under optimal taxation for the anti mgmt literature-based calibration, as pictured in Figure C.6. Under the first-best (left panel) nearly everyone, except at the very top of the ability distribution, is in Academia/Science. At the very top there is somewhat more diversity, with a significant representation especially in Management and Consulting, but still the bulk of talented individuals are allocated to Academia/Science. This is striking, but is driven by the very large positive externalities of Academia/Science the literature perceives. The large gains from the first-best allocation of talent discussed in Subsection 3.5.4 are driven by this radical reallocation.

Figure C.5: Allocation of talent by ability quantile under laissez-faire

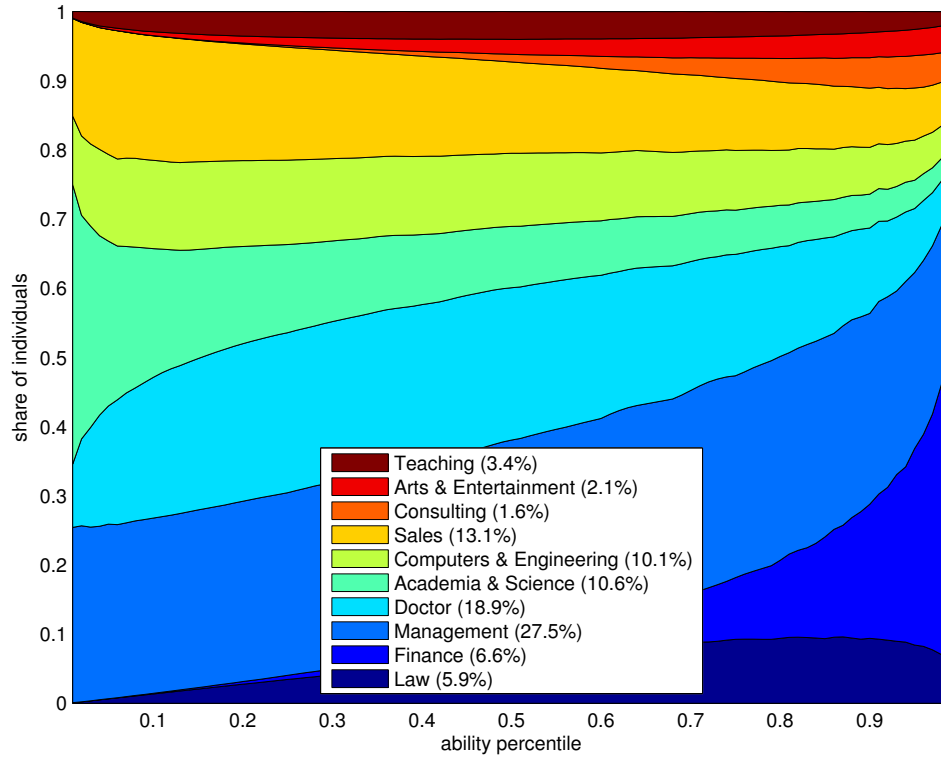
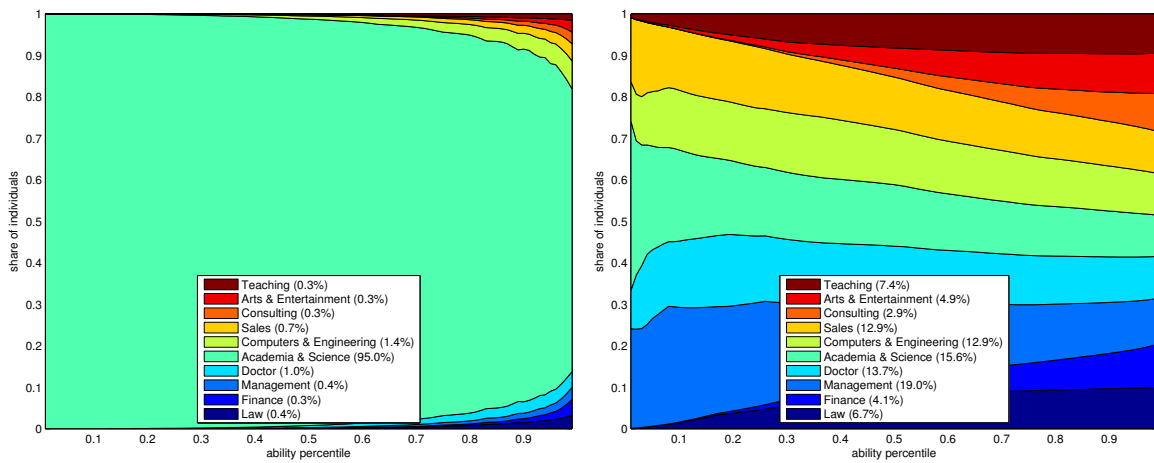
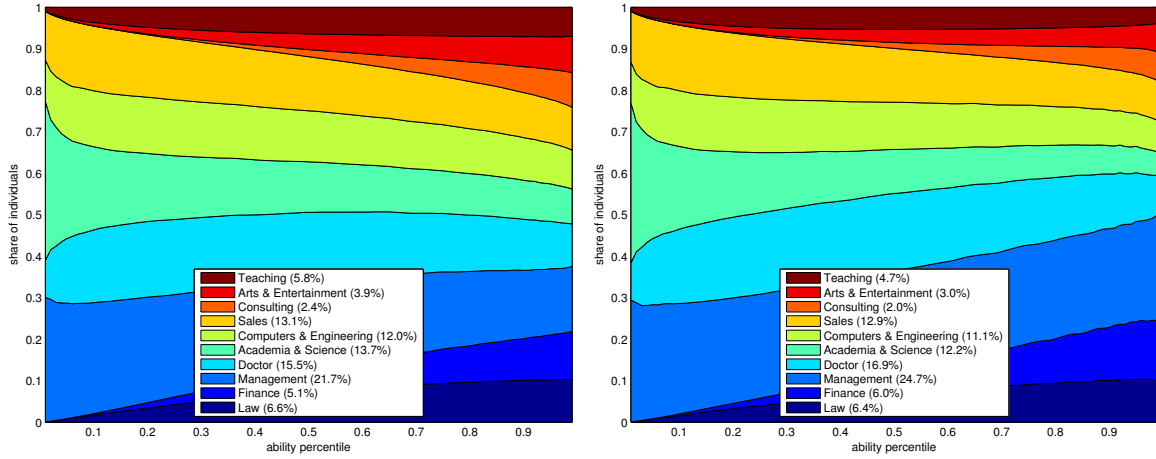


Figure C.6: Allocation of Talent in Different Regimes



Notes: Shares of individuals of various ability levels allocated to different professions under first-best, profession-specific taxation (left) and second-best, non-discriminatory taxation (right).

Figure C.7: Reallocation of Talent from Reagan Tax Reforms



Notes: Shares of individuals of various ability levels allocated to different professions under before (left) and after (right) the Reagan tax reforms.

The reallocation of talent under the second-best (right panel) is much less targeted, resulting both in a much smaller welfare gain (as discussed in Subsection 3.5.4) and a much more even distribution of talent across professions as shown in the right panel of Figure C.6. The distribution across professions at low incomes is largely unaffected, but throughout the rest of the ability distribution there is a broad heterogeneity of occupational choice. “Middle class” professions persist high up into the ability distribution, unlike under *laissez-faire*.

The contrast between the allocation of talent before and after the Reagan tax reforms, shown in Figure C.7, is a smoother and less dramatic version of the contrast between the *laissez-faire* and second-best allocations. The Reagan reforms cause Management and Finance to expand at the expense of the middle class professions especially at high ability levels, but not nearly as dramatically as under *laissez-faire*.

C.8 Proofs

Proof of Proposition 8. We begin by considering the case of pure career-switching and no intensive elasticity. First, note that the adjustment of marginal tax rates at any level must

leave average externalities at each income level constant as the average entrant and exiter has the same externality as currently prevails at that income level. Second, note that the proposed tax satisfies the necessary conditions of Proposition 7, given that there is no intensive elasticity: for any y

$$E [\Delta T + \Delta E | \partial \Theta(y)] = E \left[T(w_p h_p^*) + e_p w_p h_p^* - T(w_q h_q^*) - e_q w_q h_q^* | \partial \Theta(y) \right],$$

which equals

$$E \left[E \left[T(y'') + e_p y'' - T(y') - e_q y' | \partial \Theta(y), w_q h_q^* = y', w_p h_p^* = y'' \right] | \partial \Theta(y) \right],$$

which itself equals

$$E \left[E \left[T_0 - T_0 - E[e_{p^*} | \theta(y'')] y'' + E[e_{p^*} | \theta(y')] y' + e_p y'' - e_q y' \right. \right. \\ \left. \left. | \partial \Theta(y), w_q h_q^* = y', w_p h_p^* = y'' \right] | \partial \Theta(y) \right] = 0.$$

Next note that any other continuous tax scheme must violate these conditions. Suppose, to the contrary, that there is another scheme \hat{T} that obeys the conditions with $\hat{T} \neq T$. Then either $\hat{T} = T + k$ for some constant k or we can identify an open set $Y \subset (\underline{y}, \bar{y})$ such that $\hat{T}(y) - T(y) > \hat{T}(y') - T(y') \forall y \in Y$ and $y' \in (\underline{y}, \bar{y}) \setminus Y$. We deal with each of these cases separately:

1. $\hat{T} = T + k$: In this case, the allocation of every type to a profession is the same, but the tax scheme is either infeasible (if $k < 0$) or burns money ($k > 0$).
2. \hat{T} differs substantively: To show that the conditions of Proposition 7 are violated it is sufficient to show that any weighted sum of the conditions differs from 0. Because W is open it can be written as a countable union of n (where n possibly equals ∞) open intervals $(\underline{y}_1, \bar{y}_1), (\underline{y}_2, \bar{y}_2), \dots$. Consider

$$\sum_{i=1}^n E \left[\Delta T + \Delta E | \partial \Theta(\underline{y}_i) \right] f_s(\underline{y}_i) - E \left[\Delta T + \Delta E | \partial \Theta(\bar{y}_i) \right] f_s(\bar{y}_i) \\ = \sum_{i=1}^n \int_{\partial \Theta(\underline{y}_i)} [\Delta T(\theta) + \Delta E(\theta)] f(\theta) d\theta - \int_{\partial \Theta(\bar{y}_i)} [\Delta T(\theta) + \Delta E(\theta)] f(\theta) d\theta,$$

which equals

$$\sum_{i=1}^n \int_{\partial\Theta(\underline{y}_i) \setminus (\partial\Theta(\underline{y}_i) \cap \partial\Theta(\bar{y}_i))} [\Delta T(\theta) + \Delta E(\theta)] f(\theta) d\theta \\ - \int_{\partial\Theta(\bar{y}_i) \setminus (\partial\Theta(\underline{y}_i) \cap \partial\Theta(\bar{y}_i))} [\Delta T(\theta) + \Delta E(\theta)] f(\theta) d\theta.$$

Now note that $\partial\Theta(\underline{y}_i) \cap \partial\Theta(\bar{y}_i)$ equals

$$\left\{ \theta \in \Theta : \exists p, q \in 1, \dots, N : \left(w_p h_p^* < \underline{y}_i < \bar{y}_i < w_q h_q^* \right) \wedge \right. \\ \left. \left(w_p h_p^* - T(w_p h_p^*) + \phi(h_p^*; \phi_p) = w_q h_q^* - T(w_q h_q^*) + \phi(h_q^*; \phi_q) \right) \right\},$$

which we abbreviate by $\partial\Theta_{y_p < \underline{y}_i < \bar{y}_i < y_q}$. On the other hand

$$\partial\Theta(\underline{y}_i) \setminus (\partial\Theta(\underline{y}_i) \cap \partial\Theta(\bar{y}_i)) = \partial\Theta_{y_p < \underline{y}_i < y_q < \bar{y}_i}$$

and

$$\partial\Theta(\bar{y}_i) \setminus (\partial\Theta(\underline{y}_i) \cap \partial\Theta(\bar{y}_i)) = \partial\Theta_{\underline{y}_i < y_p < \bar{y}_i < y_q}.$$

Thus all “switches” included are from incomes within Y to outside Y . Consider any income in $y \in (\underline{y}_i, \bar{y}_i)$ for some i (which is true for any $y \in Y$) and consider

$$E \left[\hat{T}(w_p h_p^*) + e_p w_p h_p^* - \hat{T}(y) - e_q y \mid w_q h_q^* \in Y, \theta \in \partial\Theta(\underline{y}_i) \cup \partial\Theta(\bar{y}_i) \right] < \\ T(w_p h_p^*) - T(w) + E \left[e_p w_p h_p^* - e_q y \mid w_q h_q^* \in Y, \theta \in \partial\Theta(\underline{y}_i) \cup \partial\Theta(\bar{y}_i) \right] = 0,$$

where the inequality follows by the definition of Y and the equality by the fact that, as above, changing the tax scheme does not change the externality properties of T as the average externality is static as individuals switch careers. Thus

$$\sum_{i=1}^n E \left[\Delta T + \Delta E \mid \partial\Theta(\underline{y}_i) \right] f_S(\underline{y}_i) - E \left[\Delta T + \Delta E \mid \partial\Theta(\bar{y}_i) \right] f_S(\bar{y}_i) < 0$$

contradicting the necessary conditions and establishing sufficiency.

Second we consider the case of pure intensive margin elasticities. In this case Proposition 7 immediately implies, given that the absence of correlation between e_{p^*} and $\epsilon_{p^*}^h$ implies the

absence of covariance, that

$$\frac{\left([E [e_{p^*}|\Theta(y)] + T'(y)] E [\epsilon_{p^*}^h|\Theta(y)] \right) f(y)}{1 - T'(y)} = 0 \implies E [e_{p^*}|\Theta(y)] = -T'(y),$$

as $f(y), E [\epsilon_{p^*}^h|\Theta(y)] > 0$. This is also sufficient because hours decrease with taxes, establishing the necessary concavity.

□